

CaricatureGS: Exaggerating 3D Gaussian Splatting Faces with Gaussian Curvature

Anonymous 3DV submission

Paper ID 382



Figure 1. Photorealistic 3D caricature avatars produced by our method.

Abstract

A photorealistic and controllable 3D caricaturization framework for faces is introduced. We start with an intrinsic Gaussian curvature-based surface exaggeration technique, which, when coupled with texture, tends to produce over-smoothed renders. To address this, we resort to 3D Gaussian Splatting (3DGS), which has recently been shown to produce realistic free-viewpoint avatars. Given a multiview sequence, we extract a FLAME mesh, solve a curvature-weighted Poisson equation, and obtain its exaggerated form. However, directly deforming the Gaussians yields poor results, necessitating the synthesis of pseudo-ground-truth caricature images by warping each frame to its exaggerated 2D representation using local affine transformations. We then devise a training scheme that alternates real and synthesized supervision, enabling a single Gaussian collection to represent both natural and exaggerated avatars. This scheme improves fidelity, supports local edits, and allows continuous control over the intensity of the caricature. In order to achieve real-time deformations, an efficient interpolation between the original and exaggerated surfaces is introduced. We further analyze and show that it has a bounded deviation from closed-form solutions. In both quantitative and qualitative evaluations, our results outperform prior work, delivering photorealistic, geometry-controlled caricature avatars.

1. Introduction

Face caricaturization refers to the action of exaggerating distinctive facial features while preserving identity. Despite its promise for lifelike, immersive avatars, producing such exaggerations in controllable, photorealistic 3D remains an open challenge. Successful mesh-based approaches are based on geometric deformations with curvature-based methods, such as the scale-aware Poisson framework [30]. When such deformed surfaces are rendered through traditional mesh-centric pipelines, such as texture mapping, the results often appear unnatural [30]. Recently, 3D Gaussian Splatting (3DGS) [18] has emerged as a potential multiview representation that provides state-of-the-art real-time photorealism by optimizing Gaussian primitives directly from a given set of images taken from various directions. This raises the following question. *Can we combine curvature-based geometric fidelity with 3DGS to generate photorealistic caricatures?*

To address this, we start with a multiview video of a subject and its extracted FLAME mesh [23]. From this, solving the weighted Poisson equation gives us the deformed caricature mesh. We rig Gaussians to the original undeformed surface and train them following a framework previously proposed for facial expressions [21]. Later, at inference, we deform the original mesh and its rigged Gaussians according to the caricature mesh, stretching, shearing, and rotating them. However, modeling these deformations as merely an additional expression, using Gaussians optimized only on

054 the input sequence, leads to low fidelity (see Fig. 5), re-
055 vealing a domain gap in which caricatures lie outside the
056 distribution of natural expression dynamics.
057

058 To bridge this gap and in the absence of real caricature
059 training data, we synthesize pseudo-ground truth (GT^*)
060 by warping each input frame with *Local Affine Transfor-
061 mations* (LAT) induced by the correspondence from the
062 original mesh to its curvature-exaggerated counterpart, pro-
063 ducing photorealistic supervision (see Sec. 3.2). During
064 training, we stochastically alternate between real views and
065 GT^* views so that a single Gaussian set jointly models both
066 natural and caricatured deformations, allowing the Gaus-
067 sians to benefit from real ground truth while adapting to
068 GT^* . To mitigate occlusion-related artifacts and protect
069 fine structures (e.g. hair and mesh boundaries), we apply
070 a spatial mask that freezes the affected Gaussians during
071 GT^* steps (Fig. 7). These Gaussians are updated only from
072 real frames, allowing a consistent appearance to accumulate
073 in their attributes.

074 Although trained only on the two sets of views, the op-
075 timized model offers additional flexibility and control at
076 inference. First, it generalizes across a continuous range
077 of caricature intensities, with the exaggeration level con-
078 trolled by an efficient linear interpolation as an approxi-
079 mation of the solution to the weighted Poisson equation, a
080 property that we demonstrate both theoretically and empiri-
081 cally. Moreover, this representation is robust to both global
082 and local deformations, enabling controlled localized edits,
083 such as exaggerating the nose size, while leaving unrelated
084 regions unchanged.

085 The new 3DGS animatable representation is the first, to
086 our knowledge, to enable photorealistic caricature render-
087 ing while faithfully retaining identity under caricature de-
088 formations. We compare it to the current state-of-the-art dy-
089 namic facial reconstruction model [21], which consistently
090 achieves higher scores and qualitative results in terms of im-
091 age fidelity, structural consistency, and identity preservation
092 metrics.

093 Our contributions include,

- 094 • A novel 3DGS training scheme that uses GT^* generated
095 with local Affine transformations that represent real and
096 caricature avatars.
- 097 • Curvature-weighted deformation with rigged 3DGS for
098 identity-preserving photorealistic caricatures.
- 099 • Real-time avatars supporting variable exaggeration levels
and fine-grained local control of facial features.

100 2. Related Work

101 2.1. Representation for 3D Head Avatars

102 Neural implicit representations have become a dominant ap-
103 proach for high-fidelity 3D head avatars, enabling photore-
104 alistic view synthesis from sparse multiview observations.

105 IMAvatar [45] combines 3D morphable-model param-
106 eters for pose and expression control using neural blend-
107 shapes and skinning fields to produce animatable head
108 avatars. ImFace [43] disentangles identity and expression
109 using two deformation fields applied to a signed distance
110 function (SDF) template. ImFace++ [44] extends this ap-
111 proach with a two-stage refinement framework that im-
112 proves detail preservation.

113 NeRFs [24] map spatial coordinates and viewing direc-
114 tions to radiance and density and render images via volu-
115 metric integration. For head avatars, Wang et al. [37] en-
116 code sparse views into a 3D structure-aware grid of anima-
117 tion codes refined by an MLP. Gafni et al. [9] integrate a
118 low-dimensional morphable face model with a neural scene
119 representation to obtain photorealistic, controllable avatars
120 from monocular video. Gao et al. [11] employ multilevel
121 voxel fields with low-dimensional expression coefficients to
122 capture elements beyond mesh blendshapes (e.g. hair and
123 accessories). INSTA [47] accelerates dynamic NeRF by
124 embedding it around a surface representation to obtain ani-
125 matable avatars from short monocular video and Avatar-
126 MAV [39] decouples appearance from motion via motion-
127 aware neural voxel grids.

128 3D Gaussian splatting [18] represents 3D scenes as
129 anisotropic Gaussian primitives, and renders them via dif-
130 ferentiable splatting. In the context of head avatars,
131 Rig3DGS [28] reconstructed scenes in a canonical Gaus-
132 sian space and learned 3DMM-guided deformations for ef-
133 ficient and photorealistic animation, while HeadGaS [7] ex-
134 tended the representation with blendable Gaussians whose
135 attributes adapt to expression coefficients. MeGA [35] in-
136 troduced a hybrid mesh–Gaussian design, combining splats
137 with mesh geometry for high-fidelity rendering and editable
138 head avatars. GaussianAvatars [27] bound deformable 3D
139 Gaussians to a parametric face mesh via a binding inher-
140 itance strategy, and SurFhead [21] replaced the 3D Gaus-
141 sians with 2D Gaussian surfels [16], applying Jacobian
142 Blend Skinning and polar decomposition, achieving state-
143 of-the-art results in dynamic head reconstruction.

144 2.2. Mesh Deformation and Exaggeration

145 Classical mesh-based approaches realize deformations us-
146 ing geometry processing, e.g., Poisson/Laplacian editing
147 and related curvature-driven deformations [19, 32, 33, 42].
148 For faces, mesh-based deformation and caricaturization
149 have been explored through both geometry-driven and data-
150 driven approaches, evolving from early parametric face
151 models to modern neural deformation networks. Early work
152 by Blanz and Vetter [2] introduced the 3D Morphable Face
153 Model (3DMM), representing shape and texture as linear
154 combinations of example faces, enabling identity and ex-
155 pression manipulation. In the caricature domain, Bren-
156 nan [3] developed an interactive system for producing line-

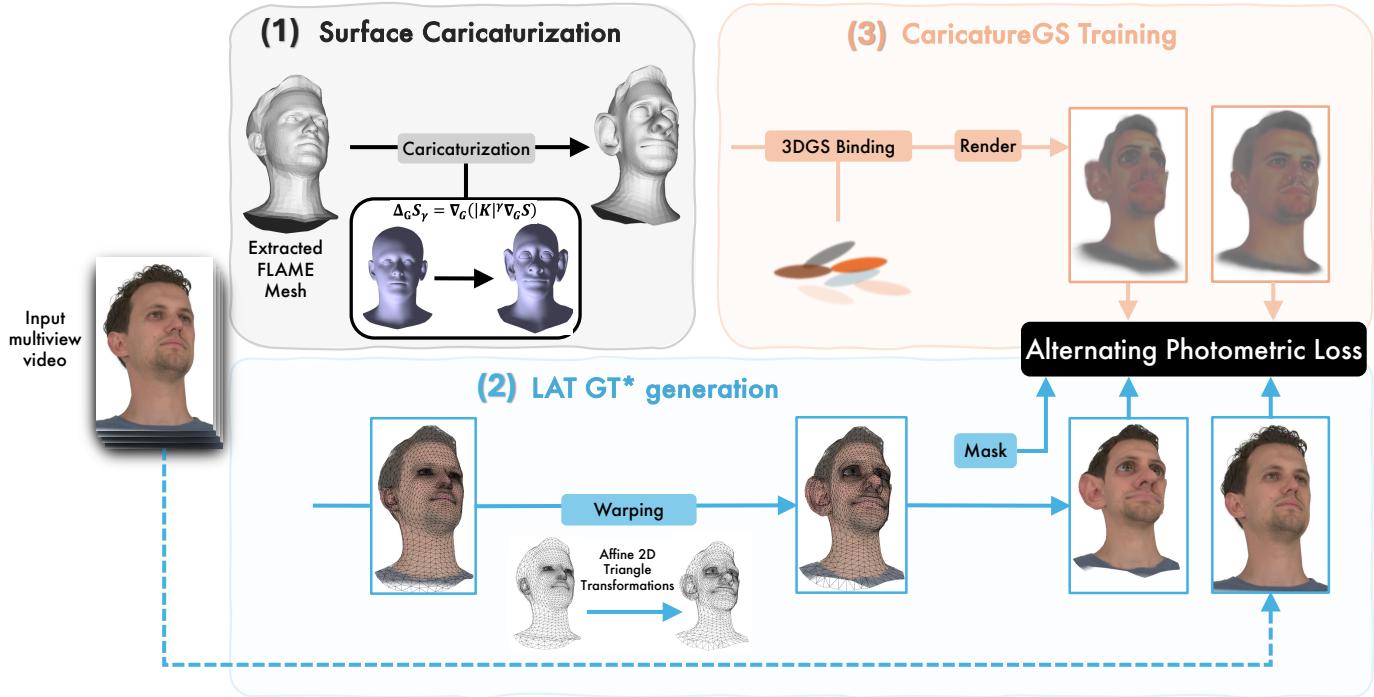


Figure 2. CaricatureGS generation framework. (1) From a subject’s multi-view video, we extract a FLAME mesh and compute a curvature-driven caricature based on it. Combined with subject-specific FLAME parameters, this yields the subject’s caricature mesh. (2) Per-triangle 2D affine transforms map the neutral mesh projection to its caricatured counterpart, warping each frame to generate pseudo-ground-truth image pairs. (3) Anisotropic 3D Gaussians primitives are bound to the original mesh and transformed to the caricature mesh via the corresponding 3D triangle transforms. Rendered neutral and caricature views are alternated and compared to their pseudo-ground-truth counterparts in joint optimization.

drawn caricatures by exaggerating the vector differences between the features of a subject and an average face. Eigen-satz [8] used curvature maps to enhance, smooth, and transfer characteristics while preserving global structure. Later, Sela et al. [31] proposed a scale-aware Poisson-based curvature framework for surface caricaturization, exaggerating geometric features while maintaining spatial and temporal coherence. Data-driven methods have enabled for more expressive and automated mesh exaggerations. Wu et al. [38] learned deformation patterns from artist-created examples to generate 3D caricatures from a single 2D portrait while preserving identity. Han et al. [12] introduced *DeepSketch2Face*, where a CNN infers and refines 3D face or caricature meshes from 2D sketches, while their later work *CaricatureShop* [13] combined vertex-wise Laplacian scaling with deep learning to produce photorealistic, personalized 2D caricatures from reconstructed 3D faces. Jung et al. [17] advanced this idea by using an MLP to map latent codes to 3D displacements, supporting controlled and diverse exaggerations. More recent approaches focus on style adaptation and broader correspondences. Yan et al. [40] presented an alignment-aware 3D face morphing framework with controller-based mapping for cross-species correspondence. Olivier et al. [25] explored GAN-based style transfer

from scans to caricatures. Yoon et al. [41] proposed *LeGO*, a one-shot method that fine-tunes a surface deformation network to replicate a target style. An additional line of work that can be adapted to facial exaggeration is the generative line, exemplified by Diffusion- and GAN-based 3DGS editors [6, 22, 36], which operate primarily on appearance while leaving the underlying geometry unchanged.

3. Method

Here, we introduce a method for creating controllable photorealistic caricaturizations of human faces with 3DGS. Our pipeline, illustrated in Fig. 2, begins with a multiview video of a subject, from which we extract a FLAME-fitted mesh. In Sec. 3.1, we describe how we deform the geometry to obtain a caricaturized mesh. To supervise 3DGS training, we generate pseudo-ground-truth caricature images (GT^*) using a 2D warping scheme (Sec. 3.2). The Gaussian primitives are then rigged to both the neutral and caricatured meshes and optimized by minimizing alternating photometric losses between their renders, the original frames, and the corresponding GT^* images (Sec. 3.3). Finally, we demonstrate that this single shared Gaussian set, although trained only on these two image domains, supports real-time ren-

181
182
183
184
185
186
187

188
189
190
191
192
193
194
195
196
197
198
199
200
201
202

203 dering across a continuous range of exaggeration levels
 204 via surface interpolation and enables region-specific edits
 205 (Sec. 3.4).

206 3.1. Surface Caricaturization

207 Starting from the temporally consistent FLAME mesh ob-
 208 tained by fitting the landmarks [34], we apply a curvature-
 209 driven deformation that exaggerates facial geometry. Since
 210 the mesh maintains consistent vertex correspondences
 211 across frames, these deformations preserve temporal coher-
 212 ence. To implement this deformation, we formulate it as a
 213 weighted Poisson equation on the surface.

214 Let $S \in \mathbb{R}^3$ be a surface with metric G and Gaussian
 215 curvature $K(p)$ for $p \in S$. For $\gamma \in [0, \gamma_f]$, we define the
 216 *weighted Poisson equation*

$$217 \Delta_G S_\gamma = \nabla_G \cdot (w(\gamma) \nabla_G S). \quad (1)$$

218 We adopt the curvature-driven deformation model intro-
 219 duced by [30], whose weights are given by $w(\gamma) = |K|^\gamma$.
 220 This gives, for each γ , the following family of Poisson equa-
 221 tions :

$$222 \Delta_G S_\gamma = \nabla_G \cdot (|K|^\gamma \nabla_G S). \quad (2)$$

223 In order to derive the deformed surface we solve the PDE
 224 by the following least-squares:

$$225 \min_{\tilde{x}} \|L\tilde{x} - b\|_A^2. \quad (3)$$

226 *L* is the discrete Laplace–Beltrami operator, defined as
 227 $L = A^{-1}W$, A is a diagonal area matrix, W is the classic
 228 cotangent weight matrix and $b = \nabla_G \cdot (|K|^\gamma \nabla_G(x))$. The
 229 weighted norm is defined as $\|F\|_A^2 = \text{trace}(F^T A F)$. We
 230 denote by S_γ the solution of the weighted Poisson equation
 231 in equation 2.

232 To accommodate open surfaces, where the Gaussian cur-
 233 vature may be ill defined on ∂S or to allow precise user-
 234 controlled exaggerations as discussed in Sec. 3.4, we im-
 235 pose boundary conditions on the selected vertices, namely:

$$236 \min_{\tilde{x} \in \mathbb{R}^n} \|L\tilde{x} - b\|_A^2 \quad \text{s.t.} \quad B\tilde{x} = x^*, \quad (4)$$

237 where $B \in \{0, 1\}^{m \times n}$ selects the rows corresponding to the
 238 set of vertices and x^* are the prescribed boundary positions.
 239 The same constrained system is solved independently for
 240 the y and z coordinates.

241 An example of the resulting mesh deformation is illus-
 242 trated in part (1) of Fig. 2.

243 3.2. GT* Generation via Local Affine Transforms

244 With these deformed surfaces, the avatar’s geometry is rep-
 245 resented in caricatured form. For photorealistic rendering,
 246 we employ mesh-rigged 3DGS, detailed in Sec. 3.3. Since

247 using 3DGS without caricature optimization yields poor re-
 248 sults (Sec. 4.2), training requires ground-truth supervision
 249 images. As real caricature images do not exist, we generate
 250 pseudo-ground truth (GT^*): photorealistic caricature
 251 images that preserve identity while ensuring multiview con-
 252 sistency.

253 One possible way to obtain such supervision is one-shot
 254 stylization (e.g., Zhou et al. [46]), which narrows the natu-
 255 ral–caricature gap using a single exemplar image. However,
 256 it fails to disentangle style from pose and identity, often
 257 transferring both instead of style alone (see supplementary).
 258 We therefore propose an alternative: Local Affine Trans-
 259 formations (LAT), illustrated in part (2) of Fig. 2.

260 LAT exploits the shared connectivity of the neutral and
 261 deformed meshes, implying a per-triangle correspondence.
 262 Consider corresponding 3D triangles $X = \{X_1, X_2, X_3\} \in \mathbb{R}^3$
 263 and $Y = \{Y_1, Y_2, Y_3\} \in \mathbb{R}^3$. Let $\pi : \mathbb{R}^3 \rightarrow \mathbb{R}^2$
 264 denote the image-plane projection, with $x_i = \pi(X_i)$ and
 265 $y_i = \pi(Y_i) \in \mathbb{R}^2$. Assuming $\{x_1, x_2, x_3\}$ are non-
 266 collinear, there exists a unique affine map,

$$\Phi(\mathbf{x}) = A\mathbf{x} + \mathbf{b}, \quad A \in \mathbb{R}^{2 \times 2}, \quad \mathbf{b} \in \mathbb{R}^2, \quad (5)$$

267 such that $\Phi(x) = y$. We then used these per-triangle 2D
 268 affine transformations to map color from the original image
 269 to the 2D projection of the deformed mesh. In practice,
 270 we apply an inverse warp from each target pixel back to the
 271 original image and use bilinear interpolation to avoid empty
 272 regions.

273 Caricature deformation can reveal regions previously
 274 self-occluded in the neutral pose or occlude regions that
 275 were visible, leaving some pixels in GT^* without valid cor-
 276 respondences. To address this, we generate 2D triangle-
 277 level mask for occluded regions. In addition, because hair
 278 strays fall outside the mesh limits and cannot be warped
 279 reliably, we add the hair boundary to the mask. The final
 280 output is pseudo-ground truth (GT^*): high-quality carica-
 281 ture images that preserve identity, ensure multiview consis-
 282 tency, and provide effective supervision for 3DGS, together
 283 with masks indicating per-pixel validity (see appendix for
 284 further details).

286 3.3. CaricatureGS Training

287 We model the avatar’s appearance photorealistically using
 288 the 3D Gaussian Splatting framework [18]. Each Gaus-
 289 sian g_i stores local attributes: position μ_i , scale s_i , rota-
 290 tion r_i , opacity σ_i , and a view-dependent color c_i . At each
 291 time frame $k \in [0, N]$, the FLAME mesh $\mathcal{M} \subset \mathbb{R}^3$ is
 292 represented by triangles $\{T_j[k]\}_{j=1}^M$, where M is the num-
 293 ber of mesh faces. To ensure spatial–temporal coherence,
 294 each Gaussian G_i is linked [27] to a specific triangle T_j by
 295 a binding index b_i , converting its local attributes to world
 296 space.

297 Building on this rigged Gaussian setup, SurfHead [21]
 298 used 2D Gaussian surfels [16], which represent surfaces as
 299 oriented planar Gaussian disks, and replaced Linear Blend
 300 Skinning (LBS) with Jacobian Blend Skinning (JBS) for
 301 Gaussians deformations, namely,

$$\Sigma_i^{1/2} = \mathbf{J}_b r_i s_i, \quad \mu'_i = \mathbf{J}_b \mu_i + T_j^x$$

302 where $\mathbf{J}_b = \exp\left(\sum_{i \in adj} v_i \log(U_i)\right) \cdot \sum_{i \in adj} v_i P_i$, (6)

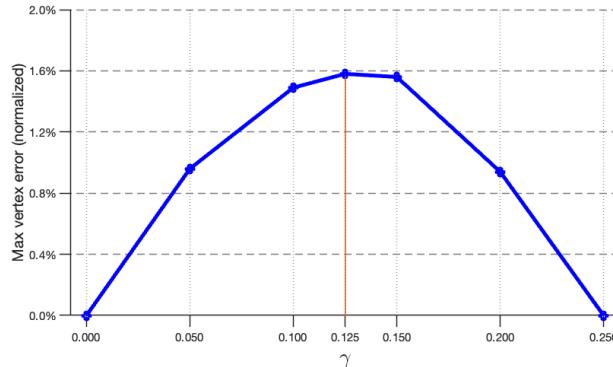
303 where v_i are learned weights and T_j^x is the triangle's
 304 barycentric center. U_i and P_i are the rotations and stretches
 305 from decomposing the Jacobian gradient \mathbf{J} via polar
 306 decomposition. Polar decomposition separates rotation and
 307 stretch, ensuring geometrically accurate Gaussian deforma-
 308 tions (see [21] for further details).

309 We show that a setup originally designed for natural
 310 facial expressions can be adapted to caricature modeling by
 311 applying the deformed caricature mesh for Gaussian de-
 312 formation and using GT* for 3DGS optimization. Never-
 313 theless, training exclusively on GT* introduces occlusion-
 314 induced artifacts and limits the model to a single ex-
 315 pression level. To overcome these limitations, we propose
 316 a joint optimization procedure that alternates supervi-
 317 sion randomly between real video frames and their carica-
 318 tured GT* counterparts, while maintaining a single shared set of
 319 Gaussians, whose rigging ensures consistent kinematics
 320 across both supervision domains. The masks introduced
 321 in Sec. 3.2 prevent supervision of Gaussians correspond-
 322 ing to caricature GT* pixels that cannot be reliably warped.
 323 The joint optimization scheme allows the caricatured 3DGS
 324 to learn beyond GT* by simultaneously filling occlusion-
 325 induced holes using supervision from the original frames.
 326 As further demonstrated in Sec. 5.2, this strategy effectively
 327 captures hair details for our caricature avatar, despite hair
 328 pixels being excluded from direct GT* supervision. More-
 329 over, as explained in Sec. 3.4, it also enables the generation
 330 of intermediate caricatures at *any* level, at inference, with-
 331 out additional capture.

332 3.4. CaricatureGS Features

333 The joint optimization not only complements the caricature
 334 Gaussians with information absent from GT* but present in
 335 the original frames, it also provides controllability advan-
 336 tages during inference.

337 **Controlling Caricature Level.** After joint training at the
 338 target exaggeration level γ_f , we empirically observe that
 339 the single-rigged Gaussian set generalizes seamlessly, ren-
 340 dering avatars from meshes deformed for any $\gamma \in [0, \gamma_f]$
 341 without additional optimization. However, obtaining the
 342 deformed mesh for each γ requires solving a curvature-
 343 weighted Poisson problem, which poses a runtime bottle-



303 Figure 3. Parametric trend of the error with respect to γ . The error,
 304 normalized by the bounding-box diagonal of the mesh, increases
 305 from both ends of γ , reaching a negligible maximum at $\frac{\gamma_f}{2}$, where
 306 $\gamma_f = 0.25$.

307 neck and makes interactive control of caricature levels im-
 308 practical. This motivates the need for a representation that
 309 can be efficiently derived from the original mesh S_0 and the
 310 precomputed caricatured mesh S_{γ_f} . We define this repre-
 311 sentation as a vertex-wise blend:

$$S_{\text{blend}}(\gamma) = (1 - \alpha) S_0 + \alpha S_{\gamma_f}, \quad \alpha \equiv \frac{\gamma}{\gamma_f}. \quad (7)$$

312 We define the residual between the approximation
 313 $S_{\text{blend}}(\gamma)$ and the exact solution $S(\gamma)$ as

$$\delta S(\gamma) = S_{\text{blend}}(\gamma) - S(\gamma). \quad (8)$$

314 In the supplementary material, we show that the L^2 energy
 315 of this residual can be bounded using Poincaré inequality
 316 together with the Lax-Milgram theorem given by

$$\|\delta S(\gamma)\|_{L^2} \lesssim \tilde{C} \gamma (\gamma_f - \gamma) \|\nabla_G S_0\|_{L^2},$$

$$\tilde{C} = C_P (\ln |K|)^2 e^{\max\{0, \gamma_f \ln |K|\}}, \quad (9)$$

317 with C_P a constant. This bound is zero at the end points
 318 $\gamma = 0, \gamma_f$, which means there is no error, as expected from
 319 (7) and maximized near $\gamma = \frac{\gamma_f}{2}$, where it remains small in
 320 practice. Empirically, we evaluate the maximal deformation
 321 error between $S_{\text{blend}}(\gamma)$ and $S(\gamma)$ on varying γ and different
 322 subjects, normalized by the mesh bounding-box diagonal.
 323 As shown in Fig. 3, the worst-case deviation is negligible,
 324 supporting the fidelity of the interpolation and confirming
 325 that it lies near the theoretical midpoint of the exaggeration,
 326 as predicted. This implies that, with this approximation, no
 327 additional Poisson equations need to be solved when infer-
 328 ring new γ values, thereby enabling full interactive control
 329 of caricature levels. In Fig. 5, we illustrate that this inter-
 330 polation scheme enables a single set of Gaussians to smoothly
 331 represent shape deformations across the full range of γ .

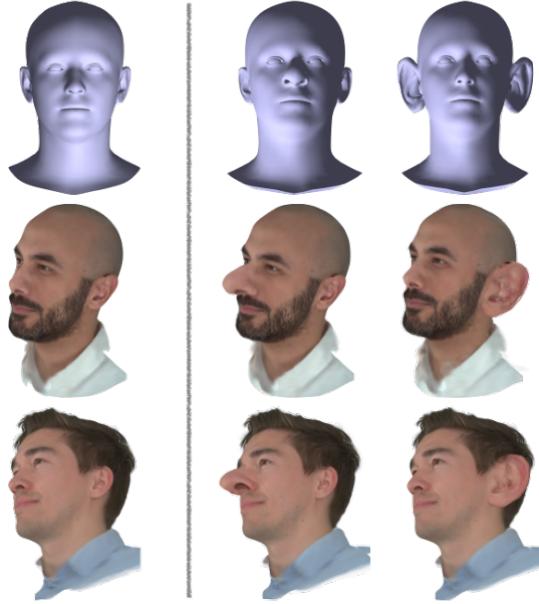


Figure 4. Visualizations of localized, semantically controlled facial exaggerations.

Localized Caricature Control. Our curvature-weighted model uses the local curvature K to generate a globally consistent caricature by solving the unconstrained Poisson equation. To target specific regions, we solve the constrained least-squares system in Eq. (4), whereby only the chosen region of interest undergoes curvature deformations, producing a smooth and localized exaggerations that blend harmonically with the rest of the face. Coupled with the training scheme in Sec. 3.3, the 3DGS, rigged to the mesh, faithfully tracks these deformations, so the same Gaussian set realizes semantically controlled exaggerations while preserving identity and global shape (see Fig. 4).

4. Experiments

We evaluate our caricaturized avatars along two main axes: (i) photorealistic rendering, (ii) identity preservation. All experiments are conducted on the NeRSemble dataset [20] and compared against the recent state-of-the-art 4D avatar reconstruction method of SurFhead [21]. Unless noted otherwise, we apply an unconstrained exaggeration with $\gamma_f = 0.25$.

4.1. Dataset

The NeRSemble dataset [20] provides a multi-view facial performance dataset captured by 16 spatially arranged, synchronized high-resolution cameras. It comprises 10 scripted sequences, 4 emotion-driven (EMO) and 6 expression-driven (EXP), plus an additional free self-reenactment sequence. For fair comparison, we adopt the same train/validation/test partition as in [21] with 120,000 training iterations. Further implementation details are provided in the supplementary.

Method	CLIP-I \uparrow	CLIP-D \uparrow	CLIP-C \uparrow	DINO \uparrow	SD \uparrow
SurFhead	0.67	0.0006	0.944	0.757	0.460
Ours	0.73	0.014	0.945	0.888	0.539

Table 1. Quantitative comparison for a caricature avatar. Higher is better for all reported metrics.

4.2. Baseline

To the best of our knowledge, there are no explicit methods that construct a dynamic 3D photorealistic model from an input multi-view video. To this end, we compare with SurFhead [21] using the authors’ official implementation. SurFhead achieves state-of-the-art performance in head reconstruction and reenactment and, in principle, can handle mesh deformations through JBS, making it the most suitable baseline for comparison. We train the SurFhead on the original input sequence and, at inference, we exaggerate the underlying mesh using γ_f , as elaborated in Sec. 2.2, thereby driving the Gaussians to represent a caricaturized avatar.

4.3. Metrics

Quantitative evaluation of caricature models is inherently challenging due to their under-constrained nature and the lack of ground-truth images. We use the following metrics for evaluation:

- **CLIP-I** (Image–Prompt Similarity) [15]: Cosine similarity between the rendered image and text in CLIP space.
- **CLIP-D** (Directional Similarity) [10]: Measures the change between source and edited images against the change between source and edited prompts.
- **CLIP-C** (Spatial Consistency): Following [14], we report CLIP image alignment between adjacent novel views of image embeddings along a novel trajectory.
- **DINO** (Identity/Structure Consistency): Following [46], we extract DINO [5] features from the renders and the corresponding original test frames and compute the cosine similarity of the embeddings.
- **SD** (Score Distillation): Inspired by DreamFusion [26], we define the reference-free metric as,

$$SD = 1 - \frac{1}{BTN} \sum_{b,t,n}^{B,T,N} \frac{\left\| \epsilon_\theta(x_t^{(b,t,n)}, t) - \epsilon_{b,t,n} \right\|_2^2}{\left\| \epsilon_{b,t,n} \right\|_2^2}. \quad (10)$$

where $\epsilon_\theta(x_t, t)$ is the noise predicted by the diffusion model [29] at time step t , ϵ is the true noise, and B, T, N refer to the image count, time step, and seed number, respectively. Higher SD indicates that the rendered image is more consistent with the training distribution of the diffusion model, which is intended to approximate the natural image distribution.

Text prompts are provided in the appendix. Together, these metrics evaluate: (i) how well the renders reflect the carica-

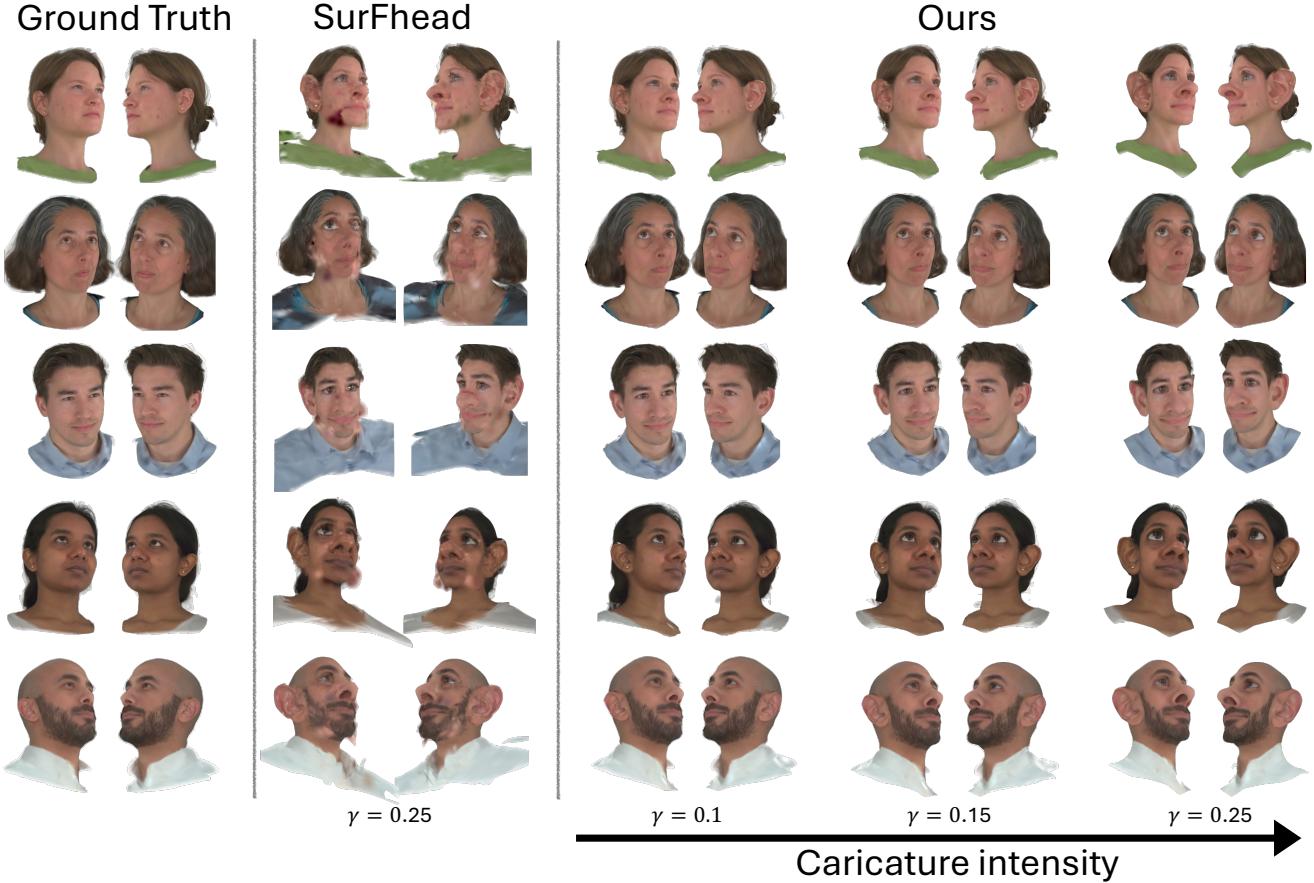


Figure 5. Rendering results from our pipeline [21]. **SURFHEAD:** Caricature generation by first reconstructing an avatar with the state-of-the-art SURFHEAD model [21], followed by mesh exaggeration. **Ours:** Renderings across different caricature intensities. Our approximation-based control interpolates smoothly along the caricature intensity axis while preserving visual fidelity.

ture intent (CLIP-I, CLIP-D, SD), (ii) identity preservation and the extent to which exaggerations remain localized to caricaturization (DINO, CLIP-D), and (iii) consistency of generated views across novel trajectories (CLIP-C).

4.4. Results

Fig. 5 presents side-by-side renderings at the target exaggeration level γ_f for our method and the baseline. Our approach maintains subject identity while delivering natural, visually pleasing exaggerations that remain consistent across views, and reduces the distortions visible in the baseline. The figure further illustrates caricature-level controllability by varying γ from 0 to γ_f , demonstrating continuous control and showing that the approximation in Sec. 3.4 successfully supports intermediate exaggeration levels.

For quantitative evaluation, we conduct a comprehensive comparison using the metrics in Sec. 4.3. As summarized in Tab. 1, our method consistently surpasses the baseline across all measures, demonstrating that the learned edits faithfully capture the intended caricature while preserving both identity and view-consistency.

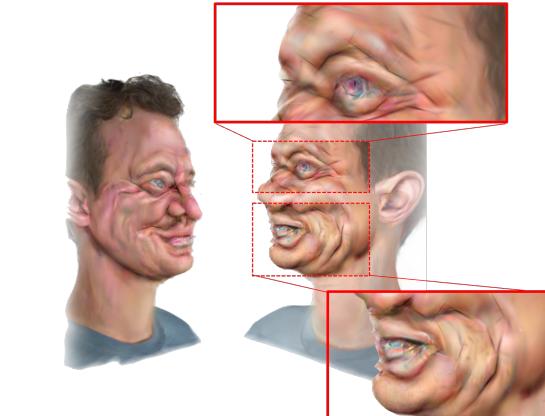
4.5. Diffusion Based Editing

As an additional baseline, we adapt a diffusion-driven, text-guided, mesh-free 3DGS editor [6] for caricaturization. Using the authors' implementation, we run 5,000 optimization steps per prompt on multiview images of a subject, guided by ControlNet-Pix2Pix. Fig. 6a presents a global edit, while Fig. 6b shows a local edit, manually masked for face and nose, respectively. While the edits appear visually plausible in individual views, it is evident that, unlike our method, this baseline suffers from (i) geometry drift, (ii) unstable, view-dependent specularities, and (iii) poor multi-view coherence.

5. Ablations

5.1. Alternated Training

In this subsection, we demonstrate that training with GT^* , generated using LAT, is essential for controlling the caricaturization level. As discussed in Sec. 4.4, training only on input images fails to generalize: rendering with a caricatured mesh yields heavily degraded outputs. In the supple-



(a) Edit instruction: “Turn him into a realistic caricature.” The result exhibits skin-tone shifts and specular degradation.



(b) Edit instruction: “Make his nose bigger.” The geometry falls apart and color inconsistencies appear across views.

Figure 6. GaussianEditor [36] caricaturization attempts. (a) Global edit. (b) Local semantic edit. Both reveal degraded geometry and appearance fidelity, particularly in novel views.

mentary, we show that training solely with GT* also fails: neutral renders appear unrealistic, with distorted Gaussian structures. These complementary failures underscore the necessity of alternating both forms of supervision for effective caricaturization control.

5.2. Mask

Due to the nature of GT* generation, certain fine details, most notably hair, are often misrepresented during the caricature stage. To address this, we identify hair regions of the mesh and freeze the corresponding Gaussian parameters with a suitable mask during GT* supervision iterations, thereby preventing updates in those regions when the caricature is rendered (see Sec. 3.2). Fig. 7 illustrates the effect: on the left, hair regions are masked and remain frozen, whereas on the right they are unfrozen and allowed to train freely, resulting in unnaturally plastic-looking hair.

6. Limitations

While our method provides a powerful framework for photorealistic 3D caricaturization, several limitations remain.

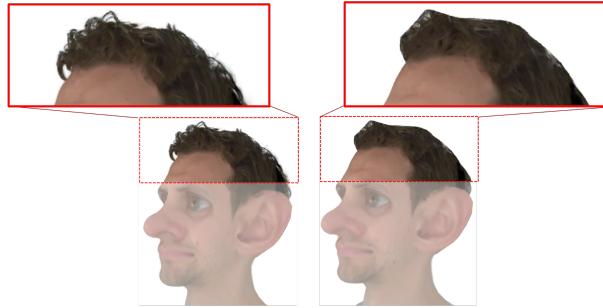


Figure 7. Ablation on hair masking. Without masking, GT* introduces visible artifacts in hair regions. Masking and freezing Gaussians associated with hair during GT* supervision effectively prevents these artifacts.

Although our approach improves upon the baseline, residual specularity artifacts persist, and small eyelid inaccuracies—amplified by over-stretching in LAT, become visually noticeable. This effect also extends to hair: training Caricature 3DGS hair with input-view supervision alone (without GT*) substantially alleviates the issue. However, in some cases, we observe slight over-smoothing of the hair. Qualitative examples of these effects are provided in the supplementary material. Finally, the deformed FLAME mesh does not fully span the space of facial expressions. For instance, eyelid closure in caricatured results is imperfect: eyes that should be completely shut under certain expressions often remain slightly open, leading to misrepresentations of eyelid geometry in the final caricature.

7. Discussion

This work demonstrates that curvature-driven geometric deformation and mesh-rigged 3D Gaussian Splatting (3DGS) can be combined into a single, controllable avatar model that remains photorealistic under large exaggerations. The key is a training scheme that alternates supervision between real views and generated pseudo-ground-truth caricature views, produced using per-triangle Local Affine Transformations (LAT) with reliability masks. One Gaussian set is capable of jointly learning both natural and caricatured appearance while retaining identity and expression. Prior work indicates that deliberate shape exaggeration can amplify discriminative geometric cues for recognition [30]. Looking ahead, we hypothesize that integrating our controllable exaggeration as a plug-in augmentation within face-recognition pipelines could improve robustness to pose and expression variability. Finally, coupling our geometry-grounded deformations with diffusion-based editors may enable semantically guided edits that are both photorealistic and extend beyond appearance-only changes to joint control of shape and appearance.

501
502
503
504
505
506
507
508
509
510
511
512
513
514

515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535

536 References

- 537 [1] Kendall E. Atkinson and Weimin Han. *Theoretical Numerical Analysis: A Functional Analysis Framework*. Springer, 538 3rd edition, 2009. 1
- 540 [2] Volker Blanz and Thomas Vetter. A Morphable Model For 541 The Synthesis Of 3D Faces. 2
- 542 [3] Susan E. Brennan. Caricature Generator: The Dynamic 543 Exaggeration of Faces by Computer. *Leonardo*, 18(3):170–178, 544 1985. Publisher: The MIT Press. 2
- 545 [4] Richard L. Burden and J. Douglas Faires. *Numerical Analysis*. Brooks/Cole, 10th edition, 2010. 1
- 547 [5] Mathieu Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, 548 Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging 549 properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on 550 Computer Vision (ICCV)*, 2021. 6
- 552 [6] Yiwen Chen, Zilong Chen, Chi Zhang, Feng Wang, Xiaofeng 553 Yang, Yikai Wang, Zhongang Cai, Lei Yang, Huaping Liu, and 554 Guosheng Lin. GaussianEditor: Swift and Controllable 555 3D Editing with Gaussian Splatting. In *2024 IEEE/CVF 556 Conference on Computer Vision and Pattern Recognition 557 (CVPR)*, pages 21476–21485, Seattle, WA, USA, 2024. 558 IEEE. 3, 7
- 559 [7] Helisa Dhamo, Yinyu Nie, Arthur Moreau, Jifei Song, 560 Richard Shaw, Yiren Zhou, and Eduardo Pérez-Pellitero. 561 HeadGaS: Real-Time Animatable Head Avatars via 3D 562 Gaussian Splatting. *arXiv e-prints*, art. arXiv:2312.02902, 563 2023. 2
- 564 [8] Michael Eigensatz, Robert W. Sumner, and Mark 565 Pauly. Curvature-Domain Shape Processing. *Computer 566 Graphics Forum*, 27(2):241–250, 2008. eprint: 567 <https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1467-8659.2008.01121.x>. 3
- 569 [9] Guy Gafni, Justus Thies, Michael Zollhöfer, and Matthias 570 Nießner. Dynamic Neural Radiance Fields for Monocular 571 4D Facial Avatar Reconstruction. *arXiv e-prints*, art. 572 arXiv:2012.03065, 2020. 2
- 573 [10] Rinon Gal, Or Patashnik, Haggai Maron, Gal Chechik, 574 and Daniel Cohen-Or. Stylegan-nada: Clip-guided 575 domain adaptation of image generators. *arXiv preprint 576 arXiv:2108.00946*, 2021. 6
- 577 [11] Xuan Gao, Chenglai Zhong, Jun Xiang, Yang Hong, Yudong 578 Guo, and Juyong Zhang. Reconstructing Personalized 579 Semantic Facial NeRF Models From Monocular Video. *arXiv 580 e-prints*, art. arXiv:2210.06108, 2022. 2
- 581 [12] Xiaoguang Han, Chang Gao, and Yizhou Yu. DeepSketch2Face: a deep learning based sketching system for 3D 582 face and caricature modeling. *ACM Transactions on Graphics*, 583 36(4):1–12, 2017. 3
- 585 [13] Xiaoguang Han, Kangcheng Hou, Dong Du, Yuda Qiu, 586 Yizhou Yu, Kun Zhou, and Shuguang Cui. CaricatureShop: 587 Personalized and Photorealistic Caricature Sketching, 2018. 588 arXiv:1807.09064 [cs]. 3
- 589 [14] Ayaan Haque, Matthew Tancik, Alexei A. Efros, Aleksander 590 Holynski, and Angjoo Kanazawa. Instruct-NeRF2NeRF: 591 Editing 3D Scenes with Instructions. *arXiv e-prints*, art. 592 arXiv:2303.12789, 2023. 6
- 593 [15] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. CLIPScore: A Reference-free Evaluation Metric for Image Captioning. *arXiv e-prints*, art. 594 arXiv:2104.08718, 2021. 6
- 595 [16] Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and 596 Shenghua Gao. 2D Gaussian Splatting for Geometrically Accurate Radiance Fields. In *Special Interest Group on Computer Graphics and Interactive Techniques Conference Conference Papers '24*, pages 1–11, 2024. arXiv:2403.17888 [cs]. 2, 5
- 597 [17] Yucheol Jung, Wonjong Jang, Soongjin Kim, Jiaolong Yang, 598 Xin Tong, and Seungyong Lee. Deep Deformable 3D Caricatures 599 with Learned Shape Control. In *Special Interest Group on 600 Computer Graphics and Interactive Techniques Conference 601 Proceedings*, pages 1–9, 2022. arXiv:2207.14593 [cs]. 3
- 603 [18] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkuehler, 604 and George Drettakis. 3D Gaussian Splatting for Real-Time 605 Radiance Field Rendering. *ACM Transactions on Graphics*, 606 42(4):1–14, 2023. 1, 2, 4
- 609 [19] ByungMoon Kim and Jarek Rossignac. Geofilter: Geometric 610 selection of mesh filter parameters. *Comput. Graph. Forum*, 611 24:295–302, 2005. 2
- 616 [20] Tobias Kirschstein, Shenhan Qian, Simon Giebenhain, Tim 617 Walter, and Matthias Nießner. NeRSemble: Multi-view 618 Radiance Field Reconstruction of Human Heads. *ACM Transactions 619 on Graphics*, 42(4):1–14, 2023. arXiv:2305.03027 [cs]. 6
- 621 [21] Jaeseong Lee, Taewoong Kang, Marcel C. Bühler, Min-Jung 622 Kim, Sungwon Hwang, Junha Hyung, Hyojin Jang, and 623 Jaegul Choo. SurFhead: Affine Rig Blending for Geometrically 624 Accurate 2D Gaussian Surfel Head Avatars, 2024. arXiv:2410.11682 version: 1. 1, 2, 5, 6, 7
- 626 [22] Guohao Li, Hongyu Yang, Yifang Men, Di Huang, Weixin 627 Li, Ruijie Yang, and Yunhong Wang. Generating Editable 628 Head Avatars with 3D Gaussian GANs, 2024. arXiv:2412.19149 [cs]. 3
- 630 [23] Tianye Li, Timo Bolkart, Michael. J. Black, Hao Li, and 631 Javier Romero. Learning a model of facial shape and 632 expression from 4D scans. *ACM Transactions on Graphics, 633 (Proc. SIGGRAPH Asia)*, 36(6):194:1–194:17, 2017. 1
- 634 [24] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, 635 Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: 636 Representing scenes as neural radiance fields for view synthesis. *CoRR*, abs/2003.08934, 2020. 2
- 638 [25] Nicolas Olivier, Glenn Kerbiriou, Ferran Argelaguet Sanz, 639 Quentin Avril, Fabien Danieau, Philippe Guillotel, Ludovic 640 Hoyet, and Franck Multon. Study on Automatic 3D Facial 641 Caricaturization: From Rules to Deep Learning. *Frontiers in 642 Virtual Reality*, 2:1–15, 2022. 3
- 643 [26] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. 644 DreamFusion: Text-to-3D using 2D Diffusion, 2022. 645 arXiv:2209.14988 [cs]. 6
- 646 [27] Shenhan Qian, Tobias Kirschstein, Liam Schoneveld, Davide 647 Davoli, Simon Giebenhain, and Matthias Nießner. GaussianAvatars: Photorealistic Head Avatars with Rigged 3D 648 Gaussians, 2024. arXiv:2312.02069 [cs]. 2, 4

- 650 [28] Alfredo Rivero, ShahRukh Athar, Zhixin Shu, and Dimitris Samaras. Rig3DGS: Creating Controllable Portraits
651 from Casual Monocular Videos. *arXiv e-prints*, art.
652 arXiv:2402.03723, 2024. 2
- 653 [29] Robin Rombach, Andreas Blattmann, Dominik Lorenz,
654 Patrick Esser, and Björn Ommer. High-resolution image
655 synthesis with latent diffusion models. In *Proceedings of
656 the IEEE/CVF Conference on Computer Vision and Pattern
657 Recognition*, pages 10684–10695, 2022. 6
- 658 [30] Matan Sela, Yonathan Aflalo, and Ron Kimmel. Computational
659 caricaturization of surfaces. *Computer Vision and
660 Image Understanding*, 141:1–17, 2015. 1, 4, 8
- 661 [31] Matan Sela, Yonathan Aflalo, and Ron Kimmel. Computational
662 caricaturization of surfaces. *Computer Vision and
663 Image Understanding*, 141:1–17, 2015. 3
- 664 [32] Olga Sorkine and Marc Alexa. As-rigid-as-possible surface
665 modeling. In *Proceedings of EUROGRAPHICS/ACM SIG-
666 GRAPH Symposium on Geometry Processing*, pages 109–
667 116, 2007. 2
- 668 [33] Olga Sorkine, Daniel Cohen-Or, Yaron Lipman, Marc Alexa,
669 Christian Rössl, and Hans-Peter Seidel. Laplacian surface
670 editing. In *Proceedings of the EUROGRAPHICS/ACM SIG-
671 GRAPH Symposium on Geometry Processing*, pages 179–
672 188. ACM Press, 2004. 2
- 673 [34] Justus Thies, Michael Zollhöfer, Marc Stamminger, Christian
674 Theobalt, and Matthias Nießner. Face2Face: Real-time
675 Face Capture and Reenactment of RGB Videos. *arXiv e-
676 prints*, art. arXiv:2007.14808, 2020. 4, 1
- 677 [35] Cong Wang, Di Kang, He-Yi Sun, Shen-Han Qian, Zi-Xuan
678 Wang, Linchao Bao, and Song-Hai Zhang. MEGA: Hybrid
679 Mesh-Gaussian Head Avatar for High-Fidelity Rendering
680 and Head Editing. *arXiv e-prints*, art. arXiv:2404.19026,
681 2024. 2
- 682 [36] Junjie Wang, Jiemin Fang, Xiaopeng Zhang, Lingxi Xie, and
683 Qi Tian. GaussianEditor: Editing 3D Gaussians Delicately
684 with Text Instructions, 2024. arXiv:2311.16037 [cs]. 3, 8
- 685 [37] Ziyan Wang, Timur Bagautdinov, Stephen Lombardi, Tomas
686 Simon, Jason Saragih, Jessica Hodgins, and Michael
687 Zollhöfer. Learning Compositional Radiance Fields of
688 Dynamic Human Heads. *arXiv e-prints*, art. arXiv:2012.09955,
689 2020. 2
- 690 [38] Qianyi Wu, Juyong Zhang, Yu-Kun Lai, Jianmin Zheng,
691 and Jianfei Cai. Alive Caricature from 2D to 3D, 2018.
692 arXiv:1803.06802 [cs]. 3
- 693 [39] Yuelang Xu, Lizhen Wang, Xiaochen Zhao, Hongwen
694 Zhang, and Yebin Liu. AvatarMAV: Fast 3D Head Avatar
695 Reconstruction Using Motion-Aware Neural Voxels. *arXiv
696 e-prints*, art. arXiv:2211.13206, 2022. 2
- 697 [40] Xirui Yan, Zhenbo Yu, Bingbing Ni, and Hang Wang. Cross-
698 Species 3D Face Morphing via Alignment-Aware Controller.
699 *Proceedings of the AAAI Conference on Artificial Intelligence*,
700 36(3):3018–3026, 2022. 3
- 701 [41] Soyeon Yoon, Kwan Yun, Kwanggyoon Seo, Sihun Cha,
702 Jung Eun Yoo, and Junyong Noh. LeGO: Leveraging a
703 Surface Deformation Network for Animatable Stylized Face
704 Generation with One Example, 2024. arXiv:2403.15227 [cs]
705 version: 1. 3
- 706 [42] Yizhou Yu, Kun Zhou, Dong Xu, Xiaohan Shi, Hujun Bao,
707 Baining Guo, and Heung-Yeung Shum. Mesh editing with
708 poisson-based gradient field manipulation. In *ACM SIG-
709 GRAPH 2004 Papers*, pages 644–651. 2004. 2
- 710 [43] Mingwu Zheng, Hongyu Yang, Di Huang, and Liming
711 Chen. ImFace: A Nonlinear 3D Morphable Face Model
712 with Implicit Neural Representations. *arXiv e-prints*, art.
713 arXiv:2203.14510, 2022. 2
- 714 [44] Mingwu Zheng, Haiyu Zhang, Hongyu Yang, Liming Chen,
715 and Di Huang. ImFace++: A Sophisticated Nonlinear 3D
716 Morphable Face Model with Implicit Neural Representations.
717 *arXiv e-prints*, art. arXiv:2312.04028, 2023. 2
- 718 [45] Yufeng Zheng, Victoria Fernández Abrevaya, Marcel C.
719 Bühler, Xu Chen, Michael J. Black, and Otmar Hilliges. I
720 M Avatar: Implicit Morphable Head Avatars from Videos.
721 *arXiv e-prints*, art. arXiv:2112.07471, 2021. 2
- 722 [46] Yang Zhou, Zichong Chen, and Hui Huang. Deformable
723 One-shot Face Stylization via DINO Semantic Guidance,
724 2024. arXiv:2403.00459 [cs]. 4, 6, 3
- 725 [47] Wojciech Zienionka, Timo Bolkart, and Justus Thies. Instant
726 Volumetric Head Avatars. *arXiv e-prints*, art.
727 arXiv:2211.12499, 2022. 2

CaricatureGS: Exaggerating 3D Gaussian Splatting Faces with Gaussian Curvature

Supplementary Material

8. Implementation considerations

Unless stated otherwise, we optimize each subject's 3D Gaussian Splatting model for 120,000 iterations, adhering to SurFhead's training protocol and evaluation split [21]. All experiments are run on a single NVIDIA RTX 3090 (24 GB VRAM). The optimization time per subject is ≈ 4 hours (this is offline training time, not rendering runtime.)

We used the NeRSembla dataset [34] with 10 subjects, 4 emotions (EMO), and 6 expressions (EXP). Expression EXP 2 is held for testing and Camera 8 serves as the validation view during training.

Caricaturization is performed once at the beginning of the training by solving the unconstrained Poisson equation, deforming the FLAME base template with $\gamma = 0.25$ (≈ 1 min).

Because FLAME uses a shared template across subjects, the deformed surface is saved and reused for all subjects. Unless stated otherwise, we report metrics over 256 frames from the rendered test sequence, aggregated across all camera viewpoints.

CLIP configuration. For text–image alignment, we use OpenAI CLIP with the ViT-B/32 backbone and the library's default preprocessing.

Prompts are: **Source:** “A realistic neutral head with natural lighting.” **Edit:** “A photorealistic caricature of a head with a highly exaggerated nose and large ears, under natural lighting.”

Defaults inherited. The optimizer, learning rate schedule, degree of spherical harmonics, and Gaussian growth/pruning follow the SurFhead [21] configuration unless otherwise specified.

9. Linear Model and Error Analysis

Notation. Let $S(u, v)$ be a parametric surface, where $(u, v) \in \mathbb{R}^2$, with a metric G and K denotes the Gaussian curvature at each point of the surface S , and

$$w(\gamma) = |K|^\gamma = e^{\gamma L}, \quad L \equiv \ln |K|. \quad (11)$$

For $\gamma \in [0, \gamma_f]$, denote by S_γ the solution of the weighted Poisson problem with Dirichlet boundary condition x^* on ∂S .

To avoid degeneracies at $K = 0$, we use ϵ to stabilize the magnitude. Note, for convenience we refer to as $|K|_\epsilon =$

$\sqrt{K^2 + \epsilon^2}$ with fixed $\epsilon > 0$. For brevity we write $|K|$ to denote this stabilized quantities.

1) Poisson equation with secant weights. The original family is defined by

$$\Delta_G S_\gamma = \nabla_G \cdot (w(\gamma) \nabla_G S). \quad (12)$$

Note, that S_0 and S_{γ_f} refer to $\gamma = 0$ and $\gamma = \gamma_f$, respectively. Define the vertex blend,

$$S_{\text{blend}}(\gamma) = (1 - \alpha) S_0 + \alpha S_{\gamma_f}, \quad \alpha \equiv \frac{\gamma}{\gamma_f}. \quad (13)$$

By linearity of Δ_G and Equation (13)

$$\begin{aligned} \Delta_G S_{\text{blend}}(\gamma) &= (1 - \alpha) \Delta_G S_0 + \alpha \Delta_G S_{\gamma_f} \\ &= \nabla_G \cdot (w_{\text{sec}}(\gamma) \nabla_G S), \end{aligned} \quad (14)$$

where *secant weight* is

$$w_{\text{sec}}(\gamma) = 1 + \frac{\gamma}{\gamma_f} (|K|^{\gamma_f} - 1). \quad (15)$$

Thus $S_{\text{blend}}(\gamma)$ solves the exact Poisson equation at level γ with $w(\gamma)$ replaced by $w_{\text{sec}}(\gamma)$, and $S_{\text{interp}}|_{\partial S} = x^*$ (see (4) for x^*).

2) Remainder and properties The secant w_{sec} is the linear interpolant of w in $[0, \gamma_f]$. By the classical interpolation remainder for C^2 functions on a closed interval (e.g., [4, Thm. 3.1], [1, §3.3]), for every $\gamma \in [0, \gamma_f]$ there exists $\xi(\gamma) \in (0, \gamma_f)$ such that

$$w_{\text{sec}}(\gamma) - w(\gamma) = \frac{w''(\xi)}{2} \gamma (\gamma_f - \gamma). \quad (16)$$

Since $w''(\gamma) = L^2 e^{\gamma L}$, we get

$$w_{\text{sec}}(\gamma) - w(\gamma) = \frac{L^2}{2} e^{\xi L} \gamma (\gamma_f - \gamma). \quad (17)$$

The secant model is exact at both endpoints (where $\alpha = 0$ and $\alpha = 1$, yielding a analytic expression in $[0, \gamma_f]$ preserving the convexity-induced non-negativity).

Since $w'' \geq 0$, $\gamma \mapsto w(\gamma)$ is convex, hence $w_{\text{sec}} - w$ is non-negative on $[0, \gamma_f]$ and vanishes at the endpoints. In particular, at $\gamma = \gamma_f/2$,

$$|w_{\text{sec}}(\frac{\gamma_f}{2}) - w(\frac{\gamma_f}{2})| \leq \frac{\gamma_f^2}{8} L^2 \max(1, e^{\gamma_f L}). \quad (18)$$

800 The maximum of this *upper bound* occurs at $\gamma_f/2$ because
 801 $\gamma(\gamma_f - \gamma)$ is maximized there.

802 3) Poincaré and Lax–Milgram for residual bound.

803 Throughout, we approximate the γ -dependent weight
 804 $w(\gamma) = |K|^\gamma$ by its secant $w_{\text{sec}}(\gamma)$ to enable a cheap vertex
 805 blend instead of solving a new Poisson problem for each
 806 γ . To justify this alternative, we should *quantify* how the
 807 weight error propagates to a *geometric residual* $\delta S(\gamma) \equiv$
 808 $S(\gamma) - S_{\text{blend}}(\gamma)$. The goal here is to derive a norm bound
 809 on δS that depends only on: (i) ellipticity and Poincaré
 810 constants of the domain, (ii) the magnitude of $\nabla_G S_0$, and (iii)
 811 the scalar secant remainder from Appendix Sec. 9. This
 812 yields a mesh and metric agnostic error budget for the blend.

813 **Setting (frozen operator).** Let (S, G) be a compact Rie-
 814 manian surface with Lipschitz boundary ∂S . We impose
 815 Dirichlet conditions $u|_{\partial S} = 0$.

816 We fix the differential operators on the surface S ,
 817 namely, the gradient and the divergence w.r.t metric G .

818 Let $V \equiv H_0^1(S)$ and define

$$\begin{aligned} a(u, v) &= \int_S \langle \nabla_G u, \nabla_G v \rangle_G dA_G \\ \|u\|_V &\equiv \|\nabla_G u\|_{L^2(S)}. \end{aligned} \quad (19)$$

819 We also define the *dual norm* by

$$\|F\|_{V'} \equiv \sup_{v \in V \setminus \{0\}} \frac{|F(v)|}{\|v\|_V}. \quad (20)$$

820 Using *Poincaré inequality*, there exists $C_P > 0$ such that,
 821 for all $u \in H_0^1(S)$,

$$\|u\|_{L^2(S)} \leq C_P \|\nabla_G u\|_{L^2(S)} = C_P \|u\|_V. \quad (21)$$

822 Hence $\|u\|_V$ is a true norm on $H_0^1(S)$ and is equivalent to
 823 the standard H^1 -norm on $H_0^1(S)$.

824 By Cauchy–Schwarz,

$$\begin{aligned} |a(u, v)| &\leq \|u\|_V \|v\|_V \quad (\text{boundedness}), \\ a(v, v) &= \|v\|_V^2 \quad (\text{coercivity with } \alpha = 1) \end{aligned} \quad (22)$$

825 where coercivity means that there exists $\alpha > 0$ such that

$$a(v, v) \geq \alpha \|v\|_V^2 \quad \forall v \in V.$$

826 **Lax–Milgram.** If a is bounded and coercive on the Hilbert
 827 space V and $F \in V'$ is bounded, then, there exists a unique
 828 solution $u \in V$, solving $a(u, v) = F(v)$ for all $v \in V$, with
 829 estimate

$$\|u\|_V \leq \frac{1}{\alpha} \|F\|_{V'} \stackrel{(22)}{=} \|F\|_{V'}. \quad (23)$$

830 For each γ , we solve the weighted Poisson PDE given by

$$\Delta_G S_\gamma = \nabla_G (w(\gamma) \nabla_G S), \quad S_\gamma|_{\partial S} = x^*. \quad (24)$$

831 Let $S_{\text{blend}}(\gamma) = (1 - \alpha)S_0 + \alpha S_{\gamma_f}$ with $\alpha = \gamma/\gamma_f$, and
 832 define

$$\begin{aligned} \psi(\gamma) &\equiv w_{\text{sec}}(\gamma) - w(\gamma) \\ \mathcal{R}_\Delta(\gamma) &\equiv \nabla_G (\psi \nabla_G S). \end{aligned} \quad (25) \quad 833$$

834 Define $F \in V'$ (weak residual functional) by

$$\begin{aligned} F(v) &= \langle \mathcal{R}_\Delta, v \rangle \\ &= \int_S (\nabla_G (\psi \nabla_G S)) v dA_G \\ &= - \int_S \psi \langle \nabla_G S, \nabla_G v \rangle_G dA_G, \end{aligned} \quad (26) \quad 835$$

836 with $v|_{\partial S} = 0$.

837 Using the dual norm and by Cauchy–Schwarz and
 838 $\|\psi\|_{L^\infty}$ -bound, we readily have

$$\begin{aligned} |F(v)| &\leq \|\psi\|_{L^\infty(S)} \|\nabla_G S\|_{L^2(S)} \|\nabla_G v\|_{L^2(S)} \\ &= \|\psi\|_{L^\infty} \|\nabla_G S\|_{L^2(S)} \|v\|_V, \end{aligned} \quad (27) \quad 839$$

840 and using (20) we get

$$\|F\|_{V'} \leq \|\psi\|_{L^\infty} \|\nabla_G S\|_{L^2(S)}. \quad (28) \quad 841$$

842 Let $\delta S \equiv S_{\text{blend}} - S_\gamma$. Subtract the weak forms for S_{blend}
 843 and S_γ to obtain

$$\begin{aligned} a(\delta S, v) &= a(S_{\text{blend}}, v) - a(S_\gamma, v) \\ &= \int_S w_{\text{sec}} \langle \nabla_G S, \nabla_G v \rangle_G dA_G \\ &\quad - \int_S w(\gamma) \langle \nabla_G S, \nabla_G v \rangle_G dA_G \\ &= \int_S \psi \langle \nabla_G S, \nabla_G v \rangle_G dA_G \\ &= - \int_S \nabla_G (\psi \nabla_G S) v dA_G \quad (*) \\ &\equiv - F(v). \end{aligned} \quad (29) \quad 844$$

845 Where in (*) we use integration by parts and Dirichlet
 846 boundary conditions on ∂S .

847 Testing with $v = \delta S$ and using coercivity and duality,

$$\begin{aligned} \|\delta S\|_V^2 &= a(\delta S, \delta S) \\ &= -F(\delta S) \leq \|F\|_{V'} \|\delta S\|_V \\ \Rightarrow \|\delta S\|_V &\leq \|F\|_{V'}. \end{aligned} \quad (30) \quad 848$$

849 Combining with the bound on $\|F\|_{V'}$ yields the *energy es-*
 850 *timate*

$$\begin{aligned} \|\delta S\|_V &\leq \|\psi\|_{L^\infty(S)} \|\nabla_G S\|_{L^2(S)} \\ \|\delta S\|_V &\leq \|w_{\text{sec}} - w\|_{L^\infty} \|\nabla_G S\|_{L^2(S)}. \end{aligned} \quad (31) \quad 851$$

859 **Optional L^2 bound.** By Poincaré on $H_0^1(S)$,

$$\begin{aligned} \|\delta S\|_{L^2(S)} &\leq C_P \|\delta S\|_V \\ &\leq C_P \|w_{\text{sec}} - w\|_{L^\infty} \|\nabla_G S\|_{L^2(S)}. \end{aligned} \quad (32)$$

861 In summary, the secant error bound yields the energy
862 bound for the residual δS by

$$\begin{aligned} \|\delta S(\gamma)\|_{L^2} &\lesssim C_P (\ln \|K\|)^2 e^{\max(0, \gamma_f \ln \|K\|)} \\ &\quad \times \gamma (\gamma_f - \gamma) \|\nabla_G S\|_{L^2(S)}. \end{aligned} \quad (33)$$

864 which depends on geometric constants of the domain (C_P).
865 The curvature in (33) is evaluated at its global maximum

$$\|K\| = K_\infty = \max_{s \in S} |K(s)| \quad (34)$$

867 We note that $S_0 = S$ (for $\gamma = 0$ by definition since there
868 is no deformation done to S), hence (33) can be written
869 using either terms.

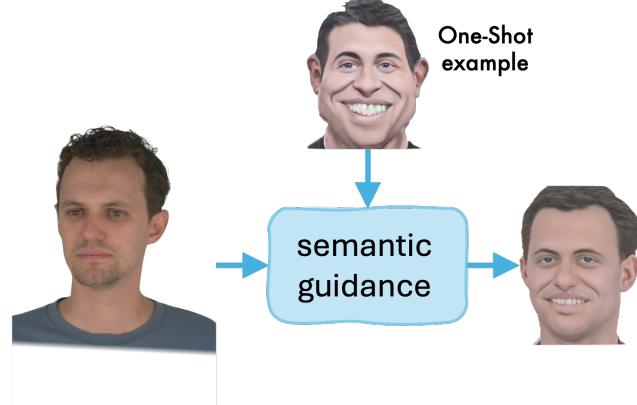
870 10. Caricature GT*via one-shot stylization

871 As discussed in Sec. 3, one-shot stylization methods (e.g.,
872 Deformable StyleGAN [46]) address the natural-caricature
873 domain gap by aligning DINO features and adapting a pre-
874 trained GAN to a single caricature exemplar. Given a target
875 style image (Fig. 8a), they synthesize stylized outputs for
876 arbitrary inputs. In practice, we observe pronounced iden-
877 tity-expression entanglement, which degrades both identity
878 fidelity and expression accuracy (Fig. 8). Moreover, the
879 outputs are not consistent across viewpoints or expressions:
880 under view changes or when transferring expressions from
881 the source, the method exhibits structural drift and a col-
882 lapsed toward the reference style (Figs. 8b and 8c), limiting
883 its suitability for our 3DGS reconstruction setting.

884 **Protocol.** We ran [46] using the official implementation,
885 employing Style1, Style2, and Style3 as target style
886 exemplars and EMO3, EMO4 for expression prompts.

887 11. Masking and GT*

888 As noted in Sec. 3.2, GT* supervision is constructed by pro-
889 jecting the FLAME mesh, fitted to each original frame, onto
890 the image. Consequently, the quality of GT* inherits any
891 mesh-image misregistration. In practice, small fitting errors
892 that are negligible at $\gamma=0$ are amplified as the caricature
893 strength increases, with the most visible drift around deli-
894 cate geometry such as the eyelids and eyeballs; see Fig. 9.
895 In addition, the deformation can reveal triangles that were
896 occluded in the original projection (e.g., along the eyelid
897 crease), creating pixels with no reliable photometric sup-
898 port.



(a) Deformable StyleGAN [46]: stylization conditioned on a target style exemplar.



(b) View variation induces identity drift and structural artifacts (e.g. neck geometry).



(c) Expressions are not preserved, outputs bias toward the style exemplar (e.g. persistent smile, forward gaze).

Figure 8. Limitations of one-shot stylization for caricature. Identity-expression entanglement and lack of view/expression consistency hinder 3DGS supervision.

To prevent these failure modes, we build a visibility-aware GT* mask. We (i) suppress supervision on triangles that become newly visible at nonzero γ relative to the original projection, and (ii) mask anatomically fragile regions prone to amplified alignment error (eyelids, ear tips).

899
900
901
902
903

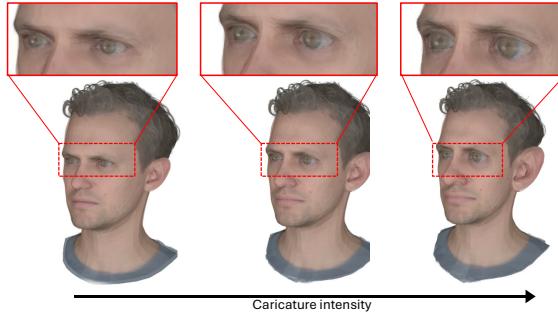


Figure 9. FLAME-image misregistration under increasing caricature strength γ . Projection drift concentrates on thin, high-curvature structures (eyelids/iris rim) and grows with γ , introducing erroneous supervision if used unfiltered.

This filtering removes inconsistent labels before they reach Gaussians anchored to those areas, yielding cleaner gradients and more stable appearance/geometry during training. The resulting GT^* thus preserves the benefits of deformation-aware supervision while avoiding artifacts introduced by projection drift and occlusions.

2. put eyelids iris break and combine it with zoom on mesh eyeballs-related to small FLAME alignment errors

12. Ablation: Alternating Supervision

Setup. As motivated in Sec. 5.1, we seek a *single* 3DGS model that renders both the original avatar ($\gamma=0$) and its caricatured counterpart ($\gamma=\gamma_f$). We compare three training schedules using identical budgets: (i) *Original-only*: supervision from original frames only. (ii) *GT^{*}-only*: supervision from caricatured (GT^*) frames only. (iii) *Alternating (ours)*: alternating mini-batches from both sources. We set the target exaggeration to $\gamma_f=0.25$ and evaluate along the interpolation path $\gamma \in \{0, 0.10, 0.15, 0.20, 0.25\}$.

Findings. Original-only (i) fits the undeformed scene well but fails to generalize to caricatured geometry Fig. 10, yielding visible distortions under nonzero γ . Conversely, GT^* -only (ii) represents the caricatured avatar but degrades markedly at $\gamma=0$. In addition, GT^* -only exhibits systematic artifacts around hair and other structures that extend beyond the tracked mesh support (*e.g.* holes or under-coverage), because those pixels are never directly supervised in the warped domain, see Fig. 11.

Our alternate schedule (iii) maintains high fidelity at both endpoints and produces smooth interpolation across γ (see Fig. 12), avoiding the hair/occlusion failures seen in (ii). Practically, alternating acts as a simple multi-domain regularizer, as it preserves appearance outside the mesh support (from original frames) while learning the exaggerated geometry and view-dependent effects required by GT^* .

Conclusions. Alternating supervision is necessary to obtain a *single* 3DGS that is faithful at $\gamma=0$ and $\gamma=\gamma_f$ and stable along the interpolation path, while training on either domain alone leads to domain-specific overfitting and characteristic failure modes.

938
939
940
941
942

3DV 2026 Submission #382. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

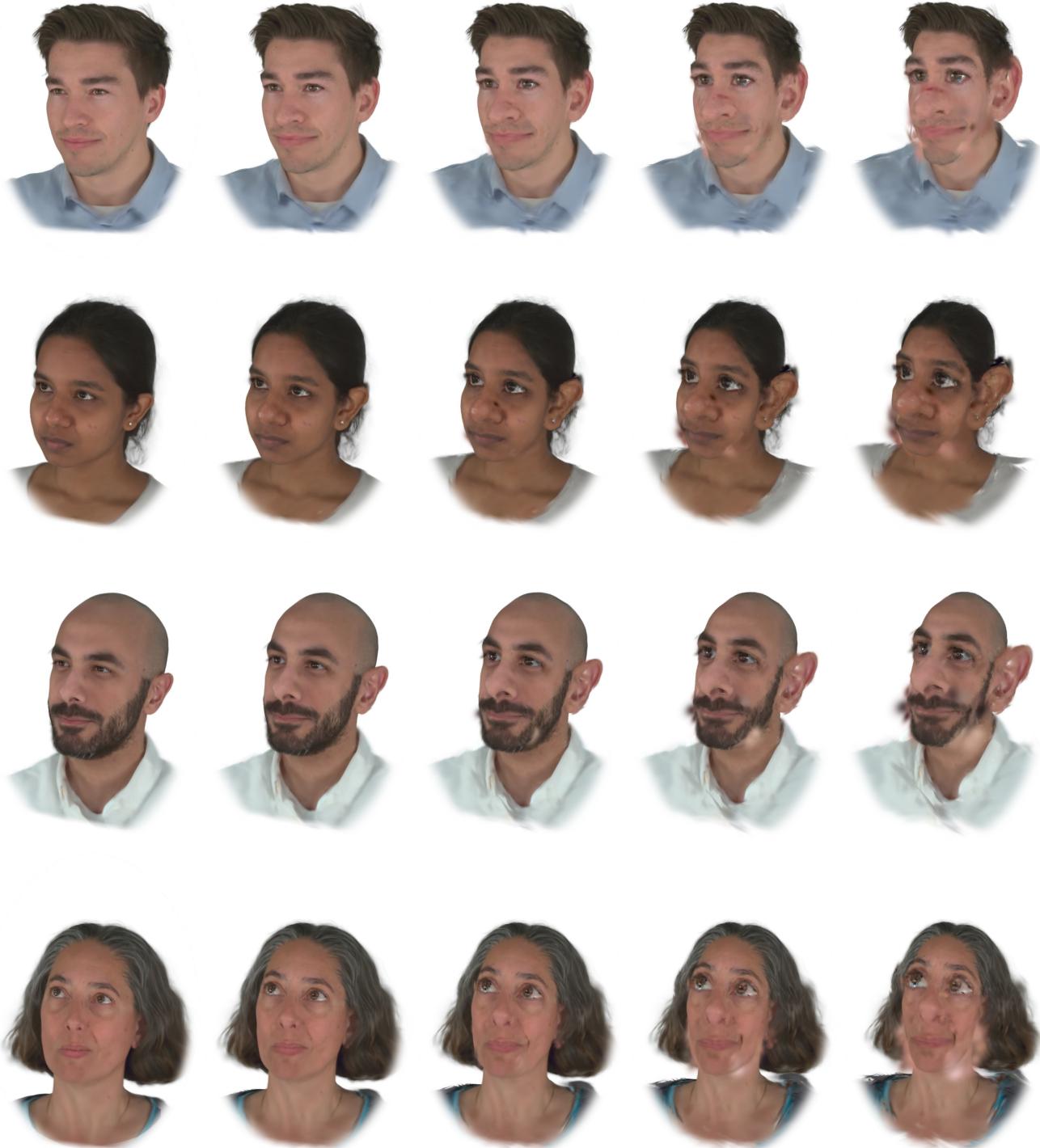


Figure 10. Training on original frames only

3DV 2026 Submission #382. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

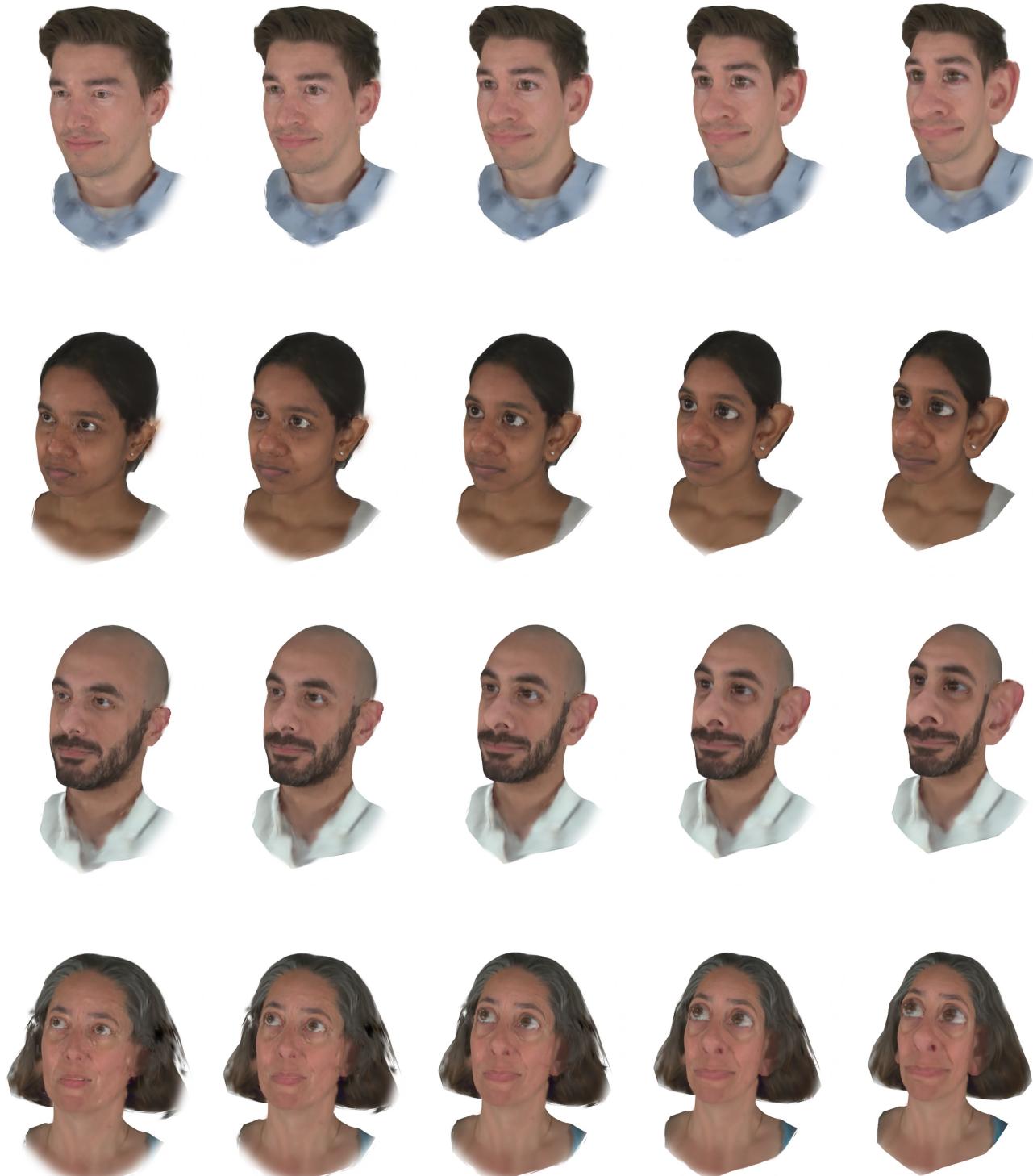


Figure 11. Training on GT*frames only.

3DV 2026 Submission #382. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

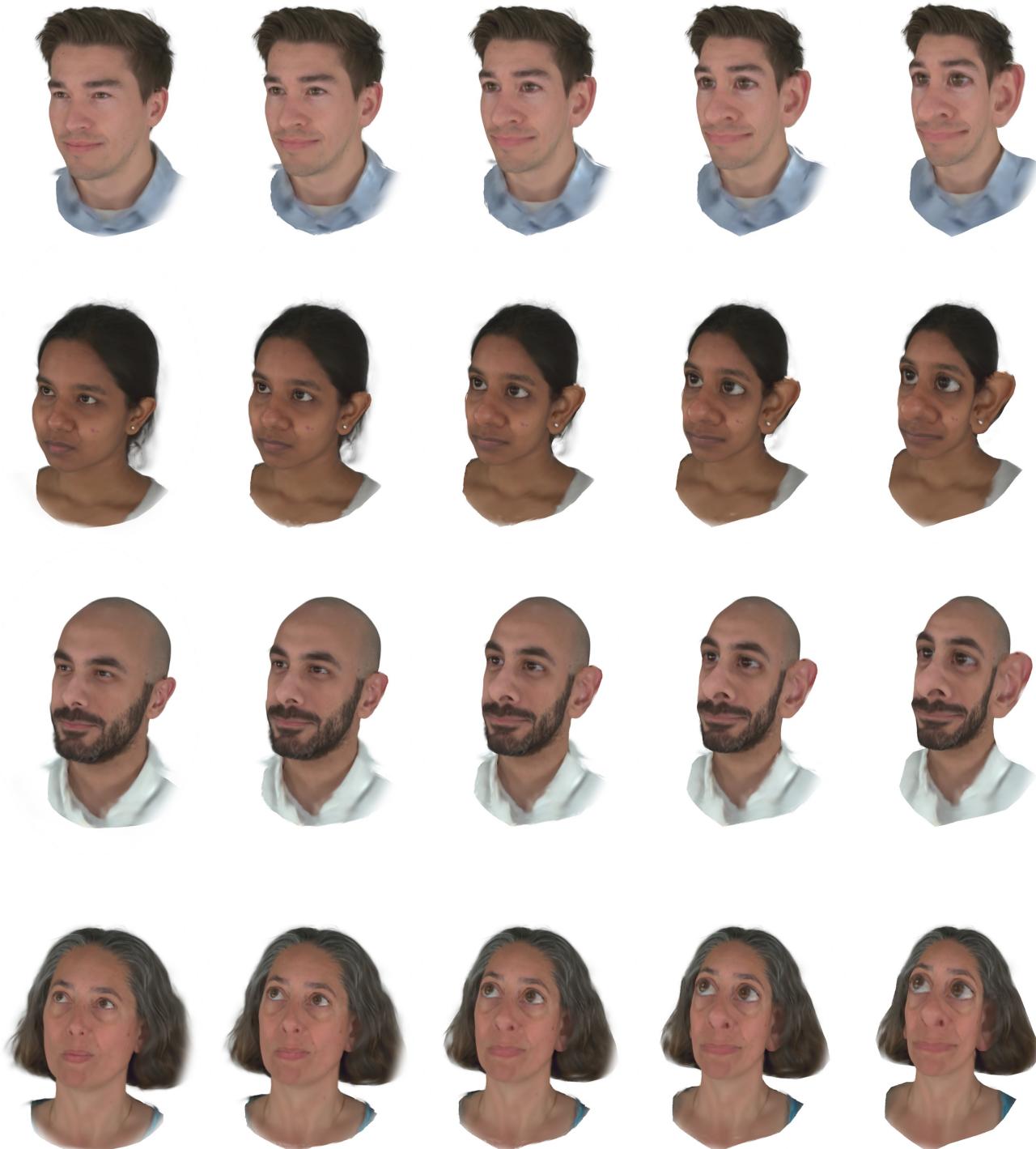


Figure 12. Training on both original and GT* frames interleaved