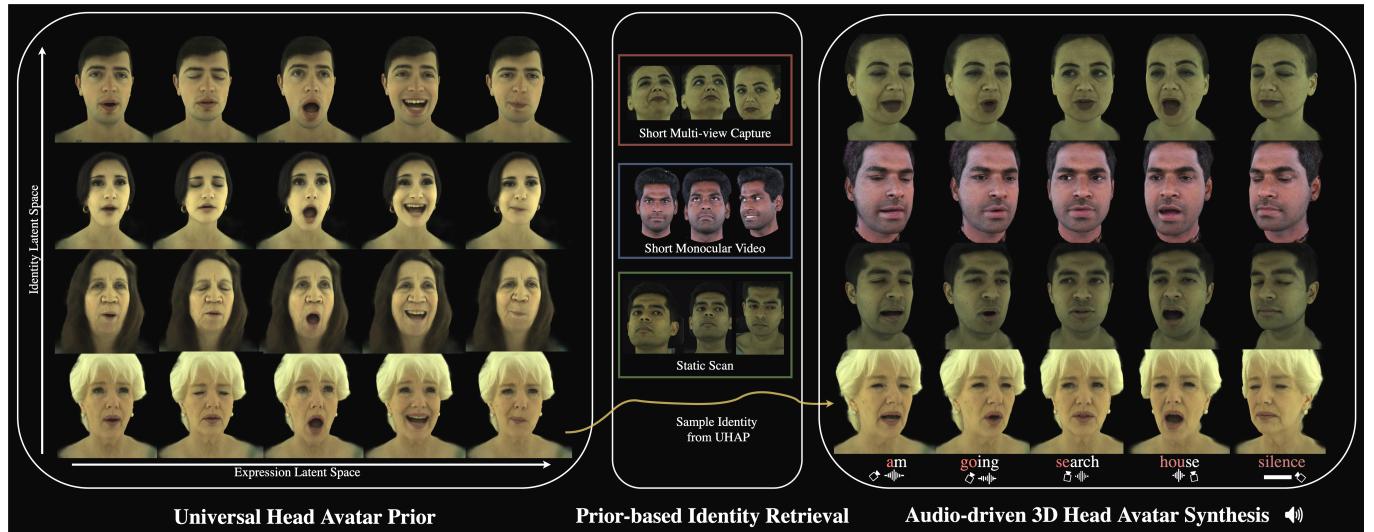


1 Audio-Driven Universal Gaussian Head Avatars

2 ANONYMOUS AUTHOR(S)

3 SUBMISSION ID: 1860



25 Fig. 1. We present a new method for audio-driven photorealistic 3D head avatar synthesis with a Universal Head Avatar Prior (**UHAP**). (Left) We learn a
26 Universal Head Avatar Prior from a diverse dataset, capturing rich facial geometry and appearance across multiple identities (rows) and dynamic expressions
27 (columns). (Center) Given minimal data for a new subject—whether a short multi-view capture, a monocular video clip, or a single static scan—we retrieve a
28 personalized identity from the **UHAP**. (Right) Conditioning the retrieved identity on an arbitrary speech waveform yields high-fidelity, lip-synced full-face
29 animations that faithfully preserve identity and expression dynamics across multiple subjects.

30 We introduce the first method for audio-driven photorealistic avatar synthesis, combining a person-agnostic speech model with our novel Universal
31 Head Avatar Prior (UHAP). UHAP is trained on cross-identity multi-view
32 videos. In particular, our UHAP is supervised with neutral scan data, en-
33 abling it to capture the identity-specific details at high fidelity. In contrast to
34 previous approaches, which predominantly map audio features to geometric
35 deformations only while ignoring audio-dependent appearance variations,
36 our universal speech model directly maps raw audio inputs into the UHAP
37 latent expression space. This expression space inherently encodes, both,
38 geometric and appearance variations. For efficient personalization to new
39 subjects, we employ a monocular encoder, which enables lightweight regres-
40 sion of dynamic expression variations across video frames. By accounting
41 for these expression-dependent changes, it enables the subsequent model
42 fine-tuning stage to focus exclusively on capturing the subject’s global ap-
43pearance and geometry. Decoding these audio-driven expression codes via
44 UHAP generates highly realistic avatars with precise lip synchronization
45 and nuanced expressive details, such as eyebrow movement, gaze shifts,
46 and realistic mouth interior appearance as well as motion. Extensive eval-
47 uations demonstrate that our method is not only the first generalizable

48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114

audio-driven avatar model that can account for detailed appearance mod-
eling and rendering, but it also outperforms competing (geometry-only)
methods across metrics measuring lip-sync accuracy, quantitative image
quality, and perceptual realism.

Additional Key Words and Phrases: Audio-driven animation, Gaussian Head Avatars

ACM Reference Format:

Anonymous Author(s). 2025. Audio-Driven Universal Gaussian Head Avatars. *ACM Trans. Graph.* 1, 1 (September 2025), 12 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnn>

1 INTRODUCTION

Synthesizing photorealistic 3D head avatars, which can be driven solely from speech, presents a compelling avenue for applications ranging from virtual communication to digital entertainment [Fan et al. 2022b]. The goal is to generate accurate lip synchronization along with expressive facial motion and, crucially, realistic visual appearance, while also ensuring temporal and view-point consistency. Achieving this using only speech as input is particularly valuable due to the lightweight sensor modality required to capture these audio signals.

A primary challenge lies in generating photorealistic appearance synchronized with accurate 3D facial motion derived from speech, while also ensuring the model generalizes to novel identities either by sampling a novel identity or by finetuning on few shot data of a real person. [Recent advancements in audio-driven video synthesis](#),

50 Permission to make digital or hard copies of all or part of this work for personal or
51 classroom use is granted without fee provided that copies are not made or distributed
52 for profit or commercial advantage and that copies bear this notice and the full citation
53 on the first page. Copyrights for components of this work owned by others than the
54 author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or
55 republish, to post on servers or to redistribute to lists, requires prior specific permission
56 and/or a fee. Request permissions from permissions@acm.org.

57 © 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
58 ACM 0730-0301/2025/9-ART
59 <https://doi.org/10.1145/nnnnnnnn.nnnnnnn>

such as VASA-1[Xu et al. 2024a], have demonstrated impressive results in generating lifelike talking faces. These methods leverage the power of diffusion-based models for mapping audio features to a video latent space to create highly realistic animations. However, these state-of-the-art approaches primarily operate in 2D, synthesizing video frames that, while visually compelling, lack the underlying 3D structure necessary for applications requiring free-viewpoint rendering. Achieving combined realism in motion and appearance in 3D remains difficult, especially without prohibitive per-person requirements like extensive multi-view capture sessions or hours of subject-specific training time [Aneja et al. 2024a; Richard et al. 2021a].

Traditional approaches to audio-driven 3D animation utilize geometric representations like 3D Morphable Models (3DMMs) [Fan et al. 2022a; Taylor et al. 2017] or artist designed template meshes [Karras et al. 2017]. While suitable for controlling basic 3D shape and motion, they face a key limitation: they do not model dynamic textures and view-dependent appearance directly from the audio signal. This deficiency makes it particularly difficult to realistically render regions such as the mouth interior or gaze shifts during speech. Consequently, the visual results often fall short of the photorealism required by many modern applications

Techniques employing modern photorealistic representations like Neural Radiance Fields (NeRFs) [Mildenhall et al. 2020] or 3D Gaussian Splatting (3DGS) [Kerbl et al. 2023a] excel at capturing appearance for static scenes or controlled dynamic captures. However, applying them directly to audio-driven animation across diverse identities often involves costly per-subject optimization or training [Aneja et al. 2024a; Ng et al. 2024; Richard et al. 2021a], requiring significant computation time (hours to days) and large amounts of per-person data, thus, hindering the creation of universal and readily deployable models. Moreover, many recent audio-driven methods, even when using the powerful diffusion models [Sun et al. 2024a; Zhao et al. 2024a], still primarily focus on driving intermediate geometry, thereby inheriting the appearance and expressiveness limitations of those representations.

Our work addresses these limitations through a novel framework centered around three key technical contributions. First, we construct a Universal Head Avatar Prior (UHAP) based on 3D Gaussian Splatting (3DGS) [Kerbl et al. 2023b]. This prior is trained on large-scale multi-view dynamic videos from studio captures, and critically incorporates supervision from neutral scan data to preserve identity-specific details during training. The resulting UHAP learns an avatar representation representation with effectively disentangled latent spaces for identity and expression. Second, unlike prior work mapping audio to intermediate geometry, we leverage a diffusion-based speech model that maps raw audio features directly into the UHAP’s expression space. A key aspect of our approach is that these predicted latent parameters explicitly encode both geometry (e.g. mouth motion) and appearance (e.g. gaze shifts) variations. Third, we enable efficient personalization of the UHAP to new subjects from sparse data, enabling practical applications such as driving a subject from a single static capture, or a short monocular video. Key to our adaptation process is a generalized monocular image encoder that estimates and factors out expression dynamics within the video frames. Thereby our monocular finetuning stage

captures the target identity’s global appearance and geometry. Importantly, our monocular image encoder does not rely on acquiring explicit geometry and appearance tracking. Decoding the audio-driven expression codes via the personalized UHAP yields the final photorealistic avatars with high-fidelity facial motion and naturally synchronized appearance changes.

In summary, our key contributions are:

- A universal framework for audio-driven and photorealistic 3D Gaussian head avatar generation allowing unconditional identity sampling as well as few-shot identity finetuning while preserving an audio-driven expression latent space.
- To this end, we first introduce a Universal Head Avatar Prior (UHAP), which effectively disentangles identity and expression latent spaces while ensuring high-fidelity synthesis thanks to our neutral texture formulation.
- A diffusion-based speech model that maps input audio to the UHAP’s expression latent space, which enables driving of the underlying 3D facial geometry and appearance. This mapping to a 3D avatar prior ensures view- and identity-consistent facial animations.
- Our monocular expression encoder facilitates a variety of few-shot identity finetuning applications such as finetuning the identity solely on a static scan or a short monocular video.

To the best of our knowledge, our work is the first that demonstrates generalization across individuals while also enabling audio-driven appearance synthesis. Our evaluation further demonstrates that we outperform geometry-only baselines in terms of audio-visual synchronization as well as visual appearance.

2 RELATED WORK

For universal avatar models to be practical for widespread adoption, they must satisfy three key criteria: they should accurately represent diverse identities, capture nuanced speech-driven expressions, and enable easy personalization from sparse observations. In what follows, we discuss prior works according to these criteria, highlighting their strengths and identifying key limitations.

2.1 Speech-driven Geometric Facial Representations

The generation of 3D facial animation from audio has a rich history, with 3D Morphable Models (3DMMs) [Blanz and Vetter 1999; Li et al. 2017] offering a generalized parametric framework for representing facial geometry and appearance. Previous approaches often involved mapping acoustic features to the parameters of these 3DMMs to achieve speech-driven animation [Aylagas et al. 2022; Daněček et al. 2022; Peng et al. 2023; Sun et al. 2024b]. However, these models are often constrained by the expressive capacity inherent in the 3DMM’s low-dimensional Principal Component Analysis (PCA) parameters, which can struggle to capture the full range of subtle, high-fidelity dynamics. Recognizing these limitations, other approaches have focused on directly modeling more detailed geometric deformations [Fan et al. 2022a; Richard et al. 2021b]. More recently, deep generative approaches, particularly diffusion models, have gained traction in this domain [Stan et al. 2023a; Sun et al. 2024b; Zhao et al. 2024b]. While powerful, many of these diffusion-based methods still

229 focus on predicting parameters for established representations like
 230 3DMMs or geometric latent models [Aneja et al. 2024b]. However,
 231 a key limitation across many speech-driven geometric represen-
 232 tations is their inability to directly model or synthesize nuanced,
 233 speech-correlated appearance changes, such as subtle gaze shifts,
 234 or deforming mouth interior. Addressing this gap, our approach
 235 synthesizes expression latents of our Universal Head Avatar Prior
 236 (UHAP), which jointly encodes, both, the subject-agnostic geometry-
 237 dependent expression changes and dynamic appearance.
 238

239 2.2 Speech-driven Appearance Methods

240 Integrating realistic, dynamic appearance with speech-driven ani-
 241 mation is crucial for photorealism but remains challenging. Early
 242 efforts primarily focused on 2D audio-driven facial animation from
 243 monocular RGB videos [Chen et al. 2018; Guan et al. 2023]. These
 244 2D methods, while achieving plausible lip sync, operate in pixel
 245 space, and, thus, they can neither achieve 3D consistency nor they
 246 support free-viewpoint rendering. Transitioning to 3D, many re-
 247 cent efforts leveraging Neural Radiance Fields (NeRF) for talking
 248 head synthesis from monocular video have shown impressive pho-
 249 torealism, such as AD-NeRF [Guo et al. 2021] and GeneFace [Ye
 250 et al. 2023], but these are often person-specific and require per-
 251 subject optimization. Other works like RAD-NeRF [Tang et al. 2022]
 252 and ER-NeRF [Li et al. 2023] focus on efficient, real-time synthesis
 253 from audio for personalized avatars. Audio-driven codec avatars,
 254 as explored in [Ng et al. 2024; Richard et al. 2021a], can produce
 255 high-fidelity personalized results but also operate on a per-subject
 256 basis. Similarly, GaussianSpeech [Aneja et al. 2024a] achieves de-
 257tailed, personalized audio-driven avatars using 3D Gaussian Splat-
 258 tting by learning expression-dependent color and dynamic wrinkles,
 259 but is tailored to individual subjects. TexTalker [Li et al. 2025b], a
 260 concurrent work to ours, generates dynamic textures aligned with
 261 speech-driven facial motion, using a high-resolution 4D dataset. It
 262 proposes a diffusion-based framework to simultaneously generate
 263 facial motions and dynamic textures from speech for personalized
 264 avatars. While TexTalker addresses dynamic textures, our work
 265 distinguishes itself by aiming for a universal prior that holistically
 266 controls, both, geometry and the broader appearance attributes cap-
 267 tured by 3D Gaussians, not limited to the tracked 2D texture maps,
 268 and allows for efficient adaptation to new individuals.
 269

270 2.3 Gaussian Avatar Representations

271 Recent works leveraging 3D Gaussian Splatting [Kerbl et al. 2023b]
 272 have introduced several powerful representations for creating per-
 273 sonalized, animatable head avatars. Foundational approaches such
 274 as GaussianAvatars [Qian et al. 2024a] rig 3D Gaussians directly to
 275 the FLAME model [Li et al. 2017], while others learn to deform a
 276 canonical set of Gaussians conditioned on global expression param-
 277 eters [Giebenhain et al. 2024; Saito et al. 2024; Teotia et al. 2024].
 278 ScaffoldAvatar [Aneja et al. 2025] achieves high fidelity rendering of
 279 faces using localized patch-based expressions. Gaussian Blendshapes
 280 [Ma et al. 2024] introduce an explicit blendshape formulation, where
 281 a full set of expression bases directly modulates Gaussian param-
 282 eters for facial animation. RGBAvatar [Li et al. 2025a] streamlines
 283 this design by predicting a reduced blendshape basis from FLAME
 284

285 expressions, yielding a more compact representation that supports
 286 efficient online training and real-time rendering. Specialized models
 287 like GaussianSpeech [Aneja et al. 2024a] animates subject-specific
 288 avatars by using a transformer model to predict audio-driven mesh
 289 deformations, which then drive the final 3D Gaussian represen-
 290 tation. While these works provide powerful representations for
 291 creating high-fidelity personalized Gaussian avatars, often requir-
 292 ing extensive per-subject data and training, our approach introduces
 293 a Universal Head Avatar Prior (UHAP). This generalizable, person-
 294 agnostic model learns a disentangled, cross-identity latent space for
 295 facial expressions that enables high-fidelity animation from mul-
 296 tiple modalities, such as an audio stream or a driving video, and
 297 supports efficient, few-shot personalization from limited input data
 298 for a new subject like a short monocular video or a static capture
 299 from multiple views.
 300

301 2.4 Universal Avatar Priors

302 Universal avatar models, capable of representing diverse identities
 303 and expressions within a unified framework, are pivotal for en-
 304 abling generalization to new individuals. Significant progress has
 305 been made in this domain. For instance, some recent universal pri-
 306 ors focus on achieving relighting capabilities alongside expressive
 307 control; URAvatar [Li et al. 2024] and VRMM [Haotian et al. 2024]
 308 are notable examples that allow for avatars to be rendered under
 309 novel illumination conditions, with VRMM also emphasizing vol-
 310 umetric representations built from data captured under controlled
 311 lighting. Other efforts, such as Authentic Volumetric Avatars by
 312 Cao et al. [Cao et al. 2022], have pushed the boundaries of creating
 313 high-fidelity, animatable volumetric avatars from inputs like phone
 314 scans. While these works show promising results in creating gen-
 315 eralizable and high-fidelity avatars, adapting them to new, unseen
 316 identities can present challenges. For example, approaches such as
 317 those by Cao et al. [Cao et al. 2022] and Li et al. [Li et al. 2024]
 318 (URAvatar) may involve extensive fine-tuning or the acquisition
 319 of dynamically tracked non-rigid facial geometry [Grassal et al.
 320 2022]. Additionally, many recent approaches based on 3D Gauss-
 321 ian Splatting map inputs to lower-dimensional expression spaces
 322 derived from, or aligned with, parametric models [Xu et al. 2024b;
 323 Zheng et al. 2024], which can constrain the overall expressiveness.
 324 Avat3r [Kirschstein et al. 2025] presents a feed-forward method
 325 for avatar creation using phone scans; however, its animation is
 326 driven by signals restricted to studio-tracked captures. Our Univer-
 327 sal Head Avatar Prior (UHAP), embedded within a 3D Gaussian
 328 Splatting framework, presents a framework for lightweight per-
 329 sonalization as well animation. It achieves lightweight personaliza-
 330 tion—facilitated by our image encoder that effectively disentangles
 331 dynamic subject-specific variations from global appearance and
 332 expression nuances—and, critically, learns a richer latent expression
 333 space directly from high-quality studio data, moving beyond the
 334 constraints of predefined parametric models. This learned expres-
 335 sion space directly modulates 3D Gaussian properties for animation
 336 from lightweight input signals such as audio or monocular images.
 337

338 339 340 341 342

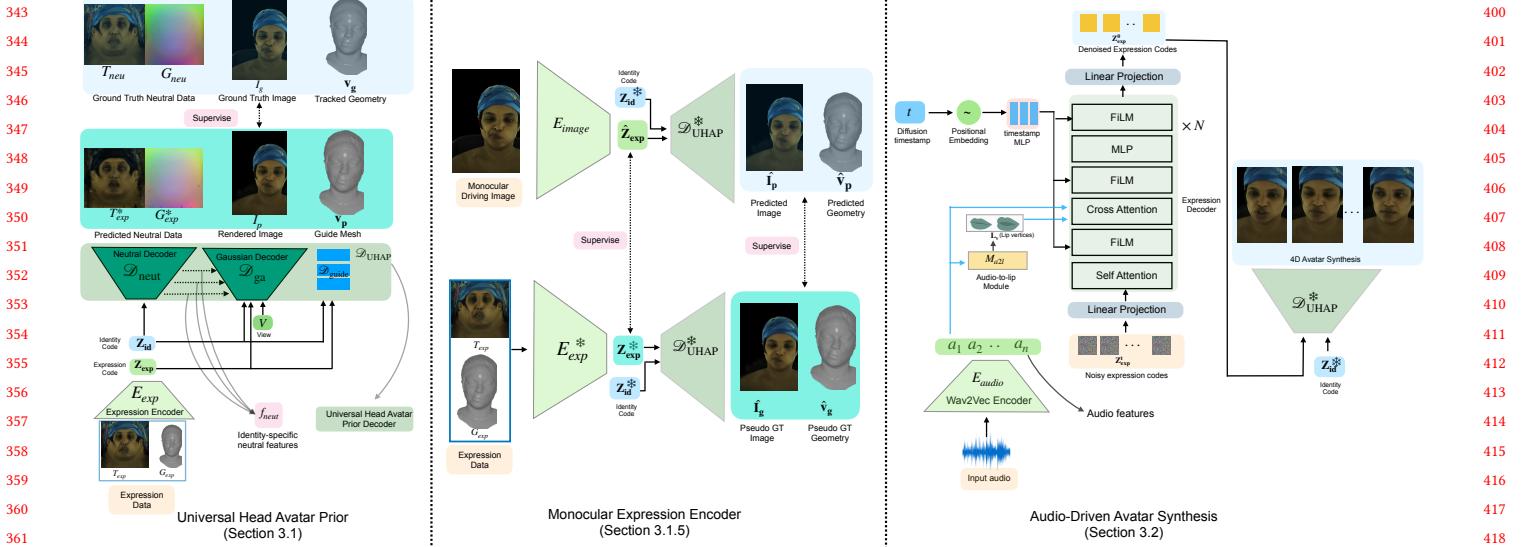


Fig. 2. Overview of our Audio-Driven Universal Gaussian Avatar pipeline. This figure illustrates the three main stages: (Sec. 3.1) Universal Head Avatar Prior (UHAP) Training: A universal decoder $\mathcal{D}_{\text{UHAP}}$ is trained on multi-identity, multi-view data to learn disentangled latent codes for identity (Z_{id}) and expression (Z_{exp}). (Sec. 3.1.5) Monocular Expression Encoder Training: An image encoder E_{image} predicts expression codes \hat{Z}_{exp} from single images, supervised by E_{exp} (from UHAP) and reconstruction losses using pseudo-ground truth data (\hat{I}_g , \hat{v}_g). (Sec. 3.2) Audio-Driven Avatar Synthesis: A diffusion model generates expression code sequences Z_{exp}^0 from audio features which, combined with Z_{id} , drive the frozen $\mathcal{D}_{\text{UHAP}}$ to synthesize the final animation.

3 METHOD

Our method synthesizes photorealistic, audio-driven 3D talking head avatars, designed for cross-identity generalization and efficient personalization. It integrates a new high-fidelity Universal Head Avatar Prior (UHAP), built upon 3D Gaussian Splatting (3DGS) [Kerbl et al. 2023b], with a universal diffusion-based speech model. This approach facilitates a direct mapping from audio features to a latent space encoding, both, expression-dependent geometry and appearance variations. Addressing the common challenge of acquiring large-scale, perfectly aligned multi-modal data (including synchronized audio, dynamic 3D geometry, and appearance across diverse identities), our framework is designed to effectively utilize varied data sources; for example, the universal prior can be trained on multi-view data even if it lacks corresponding audio, while the audio-driven aspects are learned subsequently with a diffusion model in the learned latent space. Notably, personalization is possible from sparse subject-specific video or a static scan by using a monocular expression encoding and an optimization of an identity code and decoder fine-tuning, as outlined in Fig. 2. The primary stages are: UHAP training (Sec. 3.1), learning a monocular expression encoder (Sec. 3.1.5), and audio-driven avatar synthesis (Sec. 3.2), followed by a personalization stage for new subjects (Sec. 3.3). At inference, our method requires input audio (processed into features) to drive the personalized UHAP decoder.

3.1 Universal Head Avatar Prior (UHAP)

The UHAP is our core model for synthesizing 3D Gaussian head avatars. Synthesizing realistic, drivable 3D head avatars, especially from sparse inputs or solely audio, is an inherently ill-posed problem;

UHAP addresses this by serving as a powerful, learned prior that constrains the synthesis process to enable high-fidelity and generalizable results. It is conditioned on the identity code $Z_{\text{id}} \in \mathbb{R}^{D_{\text{id}}}$ and an expression code $Z_{\text{exp}} \in \mathbb{R}^{D_{\text{exp}}}$. The identity code Z_{id} aims to capture subject-specific canonical geometry and appearance, while Z_{exp} controls facial deformations and associated appearance changes. The UHAP is trained on dense multi-view video data from multiple subjects in the Ava-256 dataset [Martinez et al. 2024], which includes registered neutral 3D scans for each subject. Unlike frameworks that encode pre-acquired neutral assets for new identities [Cao et al. 2022; Li et al. 2024], our UHAP incorporates a Neural Decoder component (Sec. 3.1.3). This network learns to inject identity-specific features—derived from the neutral scan data during UHAP training and conditioned on Z_{id} —into the main avatar decoder. This architecture promotes high-fidelity rendering. Furthermore, for new, unseen identities, it allows for efficient fine-tuning using sparse data, such as a single static scan (Sec. 3.3). Critically, our streamlined personalization strategy sidesteps an expensive and time-consuming precomputation; it does not necessitate non-rigid registration of the input dynamic or static data for these new subjects as is the case with [Cao et al. 2022; Li et al. 2024].

3.1.1 Representation. We represent each avatar as a collection of $N_g = 256k$ 3D Gaussian primitives $\{g_k\}_{k=1}^{N_g}$. Each Gaussian $g_k = \{\mathbf{t}_k \in \mathbb{R}^3, \mathbf{q}_k \in \mathbb{R}^4, \mathbf{s}_k \in \mathbb{R}^3, o_k \in \mathbb{R}_+, \mathbf{c}_k \in \mathbb{R}^{D_c}\}$ is defined by its center position \mathbf{t}_k , rotation as a unit quaternion \mathbf{q}_k , anisotropic scale \mathbf{s}_k , opacity o_k , and D_c spherical harmonics (SH) coefficients encoding the color \mathbf{c}_k . The rotation \mathbf{q}_k and scale \mathbf{s}_k together define the 3D Gaussian's covariance matrix. Images I are rendered

457 differentiably from these primitives using the Gaussian rasterizer
 458 $\mathcal{R}(\{g_k\}_{k=1}^{N_g})$, as proposed by Kerbl et al. [2023b].
 459

460 **3.1.2 Expression Encoder.** A variational autoencoder (VAE) [Stan
 461 et al. 2023b], E_{exp} , learns the expression manifold. The inputs to
 462 this encoder are UV-parameterized texture data (T) and geometry
 463 data (G). To focus on expression-specific changes, E_{exp} processes
 464 the differences: $\Delta T_{\text{exp}} = T_{\text{exp}} - T_{\text{neu}}$ and $\Delta G_{\text{exp}} = G_{\text{exp}} - G_{\text{neu}}$.
 465 These represent the deviations of the dynamic expression state
 466 ($T_{\text{exp}}, G_{\text{exp}}$) from a corresponding neutral state ($T_{\text{neu}}, G_{\text{neu}}$). The
 467 encoder maps these differences into the parameters (mean μ_{exp} and
 468 standard deviation σ_{exp}) of a multivariate Gaussian distribution:

$$\mu_{\text{exp}}, \sigma_{\text{exp}} = E_{\text{exp}}(\Delta T_{\text{exp}}, \Delta G_{\text{exp}}; \Phi_{E_{\text{exp}}}) \quad (1)$$

471 The expression code $Z_{\text{exp}} \in \mathbb{R}^{D_{\text{exp}}} (D_{\text{exp}} = 256)$ is then sampled
 472 using the reparameterization trick [Stan et al. 2023b]: $Z_{\text{exp}} = \mu_{\text{exp}} +$
 473 $\sigma_{\text{exp}} \cdot \epsilon$, where $\epsilon \sim \mathcal{N}(0, I)$.

474 **3.1.3 UHAP Decoder.** The UHAP Decoder, $\mathcal{D}_{\text{UHAP}}$, synthesizes the
 475 full 3D Gaussian avatar conditioned on the identity code Z_{id} and
 476 expression code Z_{exp} . It comprises three main components, each
 477 with its own set of learnable parameters denoted by $\Phi_{(\cdot)}$. (1) A
 478 Neutral Decoder $\mathcal{D}_{\text{neut}}$ processes Z_{id} to produce identity-specific
 479 feature maps $f_{\text{neut}} = \mathcal{D}_{\text{neut}}(Z_{\text{id}}; \Phi_{\text{neut}})$. Supervised by registered
 480 neutral 3D scan data during training, f_{neut} encapsulates the subject’s
 481 base geometry and appearance. This dedicated Neutral Decoder is
 482 a key design choice; by explicitly learning to inject these identity-
 483 specific features, it promotes better disentanglement of identity from
 484 expression, ensures more robust identity preservation during anima-
 485 tion, and contributes to more stable training, ultimately leading
 486 to sharper, higher-fidelity rendering. Critically, this enables efficient
 487 personalization to unseen captures of new identities (Sec. 3.3), with-
 488 out requiring explicit neutral 3D scans for those new subjects. (2)
 489 A Guide Mesh Decoder $\mathcal{D}_{\text{guide}}$ predicts vertex positions \hat{v}_p . Conditioned on both Z_{id} and Z_{exp} , this decoder, $\mathcal{D}_{\text{guide}}(Z_{\text{id}}, Z_{\text{exp}}; \Phi_{\text{guide}})$,
 490 predicts these as offsets relative to a canonical template mesh, v_{can} ,
 491 which has a fixed topology of 7306 vertices. (3) The Gaussian Avatar
 492 Decoder \mathcal{D}_{ga} , a CNN-based decoder, $\mathcal{D}_{\text{ga}}(Z_{\text{id}}, Z_{\text{exp}}, f_{\text{neut}}, V; \Phi_{\text{ga}})$,
 493 predicts the parameters $\{\delta t_k, q_k, s_k, o_k, c_k\}$ for the set of 3D Gaus-
 494 sians. The Gaussians are learned on a UV map that is parameter-
 495 ized by the guide mesh topology. δt_k represents predicted offsets
 496 from initial Gaussian positions which are initialized on the decoded
 497 guide mesh vertices \hat{v}_p . This decoder is conditioned on $Z_{\text{id}}, Z_{\text{exp}}$,
 498 the identity features f_{neut} (injected at various network layers), and
 499 the camera viewpoint V . The final rendered image is denoted as
 500 $I_p = \mathcal{R}(\{g_k\})$.

503 **3.1.4 UHAP Training Objective.** For UHAP training, we utilize the
 504 Ava-256 dataset [Martinez et al. 2024]. This dataset provides multi-
 505 view images for 256 subjects and, critically for our loss terms, in-
 506 cludes annotations such as: non-rigid mesh tracking for dynamic
 507 geometry (G_{exp}), which yields the ground truth vertices v_g main-
 508 taining a consistent topology across expressions and subjects; tracked
 509 dynamic appearance as UV maps (T_{exp}); and per-subject neutral
 510 scan data ($G_{\text{neu}}, T_{\text{neu}}$). The neutral data ($G_{\text{neu}}, T_{\text{neu}}$) is derived by
 511 averaging the tracked dynamic geometry and texture sequences for
 512 each subject. We jointly optimize all UHAP parameters $\Phi = (\Phi_{E_{\text{exp}}},$
 513

$\Phi_{\mathcal{D}_{\text{UHAP}}})$, where $\Phi_{\mathcal{D}_{\text{UHAP}}} = (\Phi_{\text{neut}}, \Phi_{\text{guide}}, \Phi_{\text{ga}})$. The overall loss
 514 $\mathcal{L}_{\text{UHAP}}$ is defined as following:

$$\begin{aligned} \mathcal{L}_{\text{UHAP}} = & \lambda_{\text{rec}} \mathcal{L}_{\text{rec}} + \lambda_{\text{neut}} \mathcal{L}_{\text{neut}} + \lambda_{\text{KL}} \mathcal{L}_{\text{KL}} + \lambda_{\text{geo}} \mathcal{L}_{\text{geo}} \\ & + \lambda_{\text{perc}} \mathcal{L}_{\text{perc}} + \lambda_{\text{reg_id}} \mathcal{L}_{\text{reg_id}} + \lambda_{\text{reg_gauss}} \mathcal{L}_{\text{reg_gauss}} \end{aligned} \quad (2)$$

515 Here, \mathcal{L}_{rec} is an image reconstruction loss (L1 and SSIM [Wang
 516 et al. 2004]) between the rendering I_p and ground truth image I_g .
 517 $\mathcal{L}_{\text{neut}}$ is an L1 loss on the model’s reconstruction of the neutral scan
 518 data ($G_{\text{neu}}, T_{\text{neu}}$) provided by the Ava-256 dataset. \mathcal{L}_{KL} is the KL-
 519 divergence for the VAE’s expression posterior $q(Z_{\text{exp}} | \Delta T_{\text{exp}}, \Delta G_{\text{exp}})$
 520 against $\mathcal{N}(0, I)$. \mathcal{L}_{geo} is an L2 loss comparing the predicted guide
 521 mesh vertices \hat{v}_p to the ground truth tracked vertices v_g , which are
 522 obtained from the non-rigid mesh tracking annotations in the Ava-
 523 256 dataset. $\mathcal{L}_{\text{perc}}$ is a perceptual loss [Johnson et al. 2016] between
 524 I_p and I_g . $\mathcal{L}_{\text{reg_id}}$ is an L1 norm on the identity code Z_{id} . $\mathcal{L}_{\text{reg_gauss}}$
 525 includes standard 3D Gaussian regularizations (e.g., for opacity
 526 and scale) [Teotia et al. 2024]. The $\lambda_{(\cdot)}$ values are hyperparameter
 527 weights.

528 **3.1.5 Monocular Expression Encoder.** To effectively personalize our
 529 UHAP model to new, unseen captures (videos or static images) and
 530 to facilitate image-driven animation, we train a dedicated Monocu-
 531 lar Expression Encoder, $E_{\text{image}}(I_i; \Phi_{\text{img}})$. This network’s primary
 532 role is to map an input image I_i to an estimated expression code
 533 \hat{Z}_{exp} within UHAP’s learned latent expression space. By explaining
 534 expression-dependent variations in the input image, E_{image} allows
 535 the subsequent fine-tuning of $\mathcal{D}_{\text{UHAP}}$ (Sec. 3.3) to focus on capturing
 536 the global, identity-specific attributes of the new subject.

537 E_{image} is trained using frontal images derived from our UHAP
 538 training data as inputs, with the corresponding ground truth ex-
 539 pression codes Z_{exp} (obtained from E_{exp} as described in Sec. 3.1.2)
 540 serving as targets. To enhance its generalization capabilities and
 541 encourage the learning of identity-agnostic expression features, we
 542 employ a data augmentation strategy during the training of E_{image} .
 543 This involves randomly swapping the identity of the input renderings
 544 by leveraging LivePortrait [Guo et al. 2024] as an effective
 545 expression-transfer tool to re-render the same expression on a dif-
 546 ferent identity. The objective $\mathcal{L}_{E_{\text{image}}}$ combines a squared L2 loss on
 547 the predicted latent codes with an L1 reconstruction loss. This L1
 548 loss measures the difference between renderings produced using the
 549 predicted expression code (\hat{I}_p for images, \hat{v}_p for guide mesh vertices)
 550 and pseudo-ground truth targets (\hat{I}_g, \hat{v}_g):

$$\mathcal{L}_{E_{\text{image}}} = \lambda_{\text{latent}} \|\hat{Z}_{\text{exp}} - Z_{\text{exp}}\|_2 + \lambda_{\text{recon}} (\|\hat{I}_p - \hat{I}_g\|_1 + \|\hat{v}_p - \hat{v}_g\|_1) \quad (3)$$

551 Here, \hat{I}_p and \hat{v}_p are generated using $\mathcal{D}_{\text{UHAP}}(Z_{\text{id}}, \hat{Z}_{\text{exp}})$, while \hat{I}_g and
 552 \hat{v}_g are the pseudo-ground truth targets derived from UHAP training
 553 data, as depicted in Fig. 2 (middle).

3.2 Audio-Driven Avatar Synthesis

554 To animate our Universal Head Avatar Prior (UHAP) from speech
 555 (Figure 2, Right), we generate sequences of its rich expression codes,
 556 Z_{exp} . These codes are designed to holistically modulate the entire
 557 facial state. Generating full-face expressions directly from audio,
 558 which primarily correlates with lip movements, is a key challenge.
 559 Our core audio-to-expression generator is a diffusion probabilistic
 560 model (DDPM) [Ho et al. 2020], \mathcal{G}_{θ} . For its backbone, we use the

Transformer-based model as proposed in [Ng et al. 2024]. While the framework in [Ng et al. 2024] is effective for generating expressive outputs, it was originally applied to predict person-specific codes. In contrast, our \mathcal{G}_θ is trained to synthesize sequences within our person-agnostic UHAP expression space Z_{exp} . Furthermore, our approach differs from other diffusion-based models like FaceTalk [Aneja et al. 2024b], which, though also using a Transformer architecture, primarily predicts latent codes for geometry-only parametric models. Our Z_{exp} latents, conversely, drive both the geometry and the appearance-related facial expression dynamics of UHAP. The DDPM \mathcal{G}_θ is conditioned on several inputs: audio features $A^{1:N}$, predicted lip vertices L_v , the noisy expression codes Z_{exp}^t , and the diffusion timestep t . The audio features $A^{1:N}$ are extracted from the input waveform by a Wav2Vec-based encoder [Baevski et al. 2020] (E_{audio} in Figure 2). The lip vertices L_v are predicted by a dedicated Audio-to-lip Module (M_{a2l} in Figure 2) from $A^{1:N}$ to provide strong local synchronization cues. **The Audio-to-lip module uses the Wav2Vec encoder [Baevski et al. 2020] and a pretrained, lightweight transformer to predict 338 lip vertices directly from audio. These vertices provide explicit local conditioning for our diffusion model.** The Transformer architecture [Vaswani et al. 2023] within \mathcal{G}_θ utilizes self-attention, cross-attention to fuse these conditioning signals, and FiLM layers [Perez et al. 2018] for the timestep embedding.

\mathcal{G}_θ is trained to predict the noise ϵ added to the clean expression codes Z_{exp}^0 , using the standard DDPM objective:

$$\mathcal{L}_{\text{diff}} = \mathbb{E}_{Z_{\text{exp}}, A, L_v, \epsilon, t} \left[\left\| \epsilon - \mathcal{G}_\theta \left(\sqrt{\bar{\alpha}_t} Z_{\text{exp}}^0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t, A, L_v \right) \right\|_2^2 \right] \quad (4)$$

where $\epsilon \sim \mathcal{N}(0, I)$ and $\bar{\alpha}_t$ is from the noise schedule. Training employs paired audio segments and Z_{exp} codes from the Multi-face dataset [hsin Wuu et al. 2023], where Z_{exp} codes are obtained via our UHAP’s E_{exp} (Sec. 3.1.2). During inference, the denoised sequence \hat{Z}_{exp}^0 , with a target identity code Z_{id} , drives the frozen UHAP decoder $\mathcal{D}_{\text{UHAP}}$.

3.3 Personalization for New Identities

Our UHAP can be personalized to new identities using various input data, including short dynamic captures of the new subject, or a static multi-view capture. A key advantage of our personalization approach is its efficiency and minimal data prerequisites: for the input data, we only require the rigid head pose and do not necessitate prior non-rigid 3D tracking or complex geometric registration of the subject. To adapt UHAP to a new identity, for instance from a static capture, we perform a two-stage fine-tuning process. This entire process takes approximately 20 minutes in total on a single NVIDIA A40 GPU. When adapting UHAP to a new identity from a static scan, we pass a frontal image from the scan to our pre-trained Monocular Expression Encoder E_{image} (Sec. 3.1.5) to obtain the corresponding expression code, Z_{exp} . This specific expression code Z_{exp} is then held constant throughout the subsequent two-stage fine-tuning procedure. If personalizing from a short dynamic video, E_{image} would provide per-frame expression codes. First, the identity code Z_{id} is optimized for $\sim 2k$ iterations (Fig. 12b). Second, the UHAP decoder $\mathcal{D}_{\text{UHAP}}$ is fine-tuned for $\sim 2k$ iterations (Fig. 12c) using the

fitting loss \mathcal{L}_{fit} :

$$\mathcal{L}_{\text{fit}} = \alpha_1 \mathcal{L}_{\text{photo}} + \alpha_2 \mathcal{L}_{\text{laplacian}} + \alpha_3 \mathcal{L}_{\text{offset}} + \alpha_4 \mathcal{L}_{\text{scale}} \quad (5)$$

where $\mathcal{L}_{\text{photo}}$ is an \mathcal{L}_1 photometric loss between the rendered image and the input scan; $\mathcal{L}_{\text{laplacian}}$ regularizes the smoothness of the guide mesh vertices (v_t); and $\mathcal{L}_{\text{offset}}$ and $\mathcal{L}_{\text{scale}}$ are \mathcal{L}_1 norms applied to the predicted Gaussian positional offsets and scales, respectively. The coefficients α_i are hyperparameter weights balancing these terms. This two-stage process yields the subject’s personalized neutral appearance (Fig. 12d), geometry (e), and the set of 3D Gaussians (f) that constitute the fine-tuned 3D Gaussian Avatar.

4 EXPERIMENTS

In this section, we first outline the datasets used for training our universal prior and the audio-driven synthesis model, along with key implementation details (Sec. 4.1). We then present qualitative results that demonstrate the capabilities of our method in generating audio-driven, diverse, and expressive animations (Sec. 4.2). Subsequently, we provide quantitative and qualitative comparisons against audio-driven (geometric) facial animation methods (Sec. 4.3). Finally, we conduct ablation studies (Sec. 4.4) to validate the impact of our key components and design choices within our proposed framework, such as the role of neutral features, the pretraining of our monocular encoder, and the amount of data needed for personalization.

4.1 Datasets and Implementation Details

Training Data. Our Universal Head Avatar Prior (UHAP) is trained using 230 distinct identities from the Ava-256 dataset [Martinez et al. 2024]. This dataset provides multi-view dynamic video recordings and registered neutral 3D scans for each subject *but it contains no audio*. The large number of identities in Ava-256 is crucial for learning a robust and generalizable prior (UHAP) over identity and expression. For training the audio-to-expression synthesis model, we utilize the Multiface dataset [hsin Wuu et al. 2023]. Multiface provides synchronized multi-view video data of subjects uttering a combined 650 sentences, offering rich audio-visual correspondence, though with a more limited number of identities compared to Ava-256. All identities from Multiface are unseen during the training of our UHAP prior. Multiface also provides tracked dynamic geometry (G_{exp}), dynamic appearance UV maps (T_{exp}), and corresponding neutral data (G_{neu} , T_{neu}) for its subjects. We leverage our pre-trained UHAP and its associated expression encoder (Sec. 3.1.2) to process these T_{exp} , G_{exp} sequences from the Multiface dataset, mapping them into our subject-agnostic expression space to obtain Z_{exp} codes. This allows us to create the synchronized audio-feature-to- Z_{exp} pairs necessary for training our audio-driven model. Data from 10 identities from Multiface dataset are used for training this audio model. Three Multiface dataset identities, entirely unseen by both the UHAP model and the audio model during their respective training phases, are held out exclusively for testing and evaluating the audio-visual performance of our complete pipeline.

Baseline Setup. Our method’s capability to personalize to new subjects is versatile, accommodating inputs such as static captures or dynamic videos (as detailed in Sec. 3.3). For the quantitative and qualitative comparisons against state-of-the-art methods requiring personalization (Sec. 4.3), we use a consistent setup for, both, our

Table 1. Quantitative comparison with SOTA audio-driven avatar methods on the held-out audio and universal model test subjects. Metrics are averaged across all frames and identities.

Method	PSNR \uparrow	LPIPS \downarrow	SSIM \uparrow	LSE-D \downarrow
CodeTalker	26.23	0.37	0.6518	8.30
FaceFormer	25.93	0.38	0.6475	9.32
FaceDiffuser	26.32	0.43	0.6832	8.88
Ours	27.37	0.29	0.7293	6.32

model and the baselines. Specifically, for new identities from the Multiface test set, our UHAP model is fine-tuned using approximately 500 frames (or 5 sentences) per camera from 12 views. We highlight that there is no prior work that is generalizable from speech input and enables photorealistic renderings. Instead, we compare against state-of-the-art geometry-based methods, i.e., Faceformer [Fan et al. 2022a], CodeTalker [Xing et al. 2023], and FaceDiffuser [Stan et al. 2023a]. Their output consists of animated mesh sequences in FLAME topology, which we further augmented for photorealistic rendering for fair comparison. This is achieved by training person-specific GaussianAvatars [Qian et al. 2024b] for each test identity. Crucially, these GaussianAvatars are also trained using the identical data setup as our personalization stage leverages. This ensures a fair and direct comparison in terms of the input data provided for achieving photorealistic results.

4.2 Qualitative Results

We first show the general qualitative performance of our audio-driven avatar synthesis method. Fig. 4 demonstrates the capability of our method to generate expressive audio-driven animations for three distinct synthesized identities. The sequences highlight accurate lip synchronization corresponding to the provided audio prompts, accompanied by natural-looking facial dynamics and varied expressions indicative of the spoken content. Fig. 9 illustrates the versatility and effectiveness of our personalization process across various input conditions. We show adaptation to new subjects from: a monocular video capture from the INSTA dataset [Zielonka et al. 2023] (top row), multiview captures (8 input views) from Multiface Dataset [hsin Wuu et al. 2023] (middle row), and (bottom row). We highlight that all these datasets are not part of the UHAP training. In each case, the input data (leftmost column) is used to personalize UHAP, and the subsequent audio-driven animations (right columns) demonstrate that the unique identities are well-captured and then faithfully animated with coherent speech motions. Beyond audio-driven synthesis, Fig. 11 underscores the versatility of our learned expression space through image-driven animation. Here, expressions from a source driving sequence (top row) are successfully transferred to multiple distinct target identities (rows below), demonstrating accurate expression re-targeting while consistently maintaining the unique appearance and characteristics of each target avatar. This highlights the successful latent disentanglement of the identity and expression latent space.

4.3 Comparisons with State-of-the-Art Methods

We conduct a comparative evaluation of our method against several recent state-of-the-art audio-driven facial animation techniques: Faceformer [Fan et al. 2022a], CodeTalker [Xing et al. 2023], and FaceDiffuser [Stan et al. 2023a]. As these methods primarily focus on generating 3D mesh deformations, their outputs are rendered using personalized GaussianAvatars [Qian et al. 2024b] to enable a fair photorealistic comparison. [We evaluate our method on held-out speakers/subjects from the Multiface dataset \[hsin Wuu et al. 2023\]](#).

Quantitative Comparison. Tab. 1 summarizes the quantitative results on the held-out Multiface test identities. Our method achieves superior performance across standard image reconstruction metrics, including higher PSNR and lower L1 and LPIPS [Zhang et al. 2018] scores, which indicates enhanced image fidelity and perceptual quality. Furthermore, our approach demonstrates improved audio-visual synchronization, as reflected by a better (lower) LSE-D score [Chung and Zisserman 2016]. Since these metrics test the end-to-end performance from audio to image quality, these results confirm the combined benefits of our contributions that more directly link audio input and avatar rendering.

Qualitative Comparison. Fig. 5 provides a side-by-side visual comparison against state-of-the-art methods, personalized on the same Multiface test subjects, and ground truth for specified audio segments, shown from a held-out novel viewpoint. Our method consistently produces results with higher fidelity details in terms of appearance and geometry. For instance, in subjects with facial hair, our approach renders a sharper beard that deforms naturally and coherently with speech-induced jaw and cheek movements, a detail which prior method can typically not preserve resulting in smoothed out renderings that lack photorealism. Furthermore, our model generates a significantly sharper and more realistic mouth interior, contributing to more natural expressions during speech. This, combined with more precise mouth articulation (e.g., for words like “change” and “bride”) and subtle eye movements, leads to better visual lip synchronization and overall fidelity to the ground truth, which is visibly higher compared to the baselines. This visual superiority can be attributed to our model’s ability to directly synthesize these fine-grained appearance attributes coherently with geometric deformations, all driven by the audio-driven latent expression codes.

4.4 Ablation Studies

To validate the contributions of individual components and design choices within our framework, we perform several ablation studies.

Impact of Neutral Features in UHAP. Fig. 10 evaluates the importance of incorporating identity-specific neutral features (f_{neut}) during UHAP training. The visual comparison shows renderings with our full model, without neutral features, and the Ground Truth, alongside quantitative metrics. Removing these neutral feature inputs results in a discernible degradation in rendering quality and the precision of identity preservation. This highlights the critical role of these learned neutral characteristics in achieving high-fidelity personalization with our UHAP.

Role of Pretrained Monocular Expression Encoder. The significance of employing a pretrained monocular expression encoder

(E_{image}), as opposed to training it from scratch during subject-specific fine-tuning, is demonstrated in Fig. 8. The left panel shows results when the encoder is trained during fitting, the center panel shows our approach with a pretrained encoder, and the right panel shows ground truth. When the expression encoder is trained concurrently with subject fine-tuning, it tends to learn a mapping that entangles expression with the specific subject’s appearance and geometry. This causes a mismatch when expression codes from our universal audio model, which expects the original, disentangled latent space semantics, are fed into this subject-adapted decoder, leading to distorted expressions and incorrect appearance. Our proposed approach, which utilizes the pretrained encoder designed to isolate true expression variations, maintains compatibility with the audio model’s output, ensuring faithful synthesis.

5 CONCLUSION

We have introduced a novel framework for the audio-driven synthesis of universal, photorealistic 3D Gaussian head avatars. Our Universal Head Avatar Prior (UHAP), learns a rich expression latent space that holistically controls both detailed geometry and dynamic appearance. Combined with an efficient personalization strategy adaptable to sparse inputs and an audio-to-expression diffusion model, our approach generates high-fidelity animations. These animations demonstrate accurate lip synchronization and nuanced facial dynamics such as eye-gaze shifts, all generalizing across diverse identities and sparse, monocular capture settings.

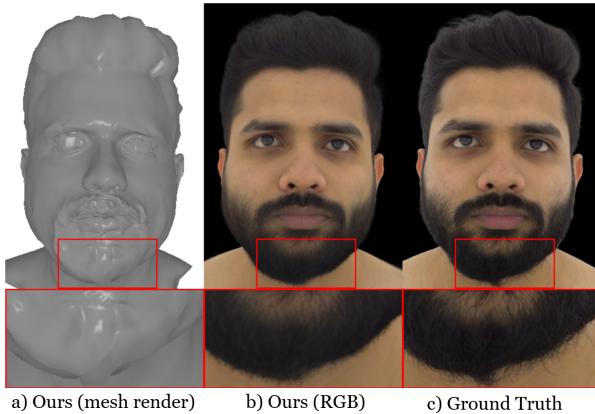


Fig. 3. Our method struggles with fine regions such as beards, where our model’s geometry-fitting (a) averages out fine details, leading to less accurate reproduction (b) compared to ground truth (c).

Limitations. Despite promising results, our method has limitations. While synthesized upper-head expressions generally align well with speech-driven mouth motion, the current gaze behavior can sometimes appear unnatural, potentially reflecting the script-reading nature of the audio training data. **Training on conversational audio-visual data and adding explicit gaze control are potential avenues for future works to overcome this limitation.** Furthermore, although our model reconstructs fine details like static hair strands, it **struggles with elements that exhibit complex, independent motion relative to the skin surface, such as beards as shown in Fig. 3, which the**

current representation may not perfectly register.

This limitation can be overcome by leveraging strand-based representation for facial hair [Winberg et al. 2022]. While our appearance model can render avatars in real-time (50 FPS on a single NVIDIA A40 series GPU), the audio-to-latent module does not decode expression codes in real-time due to the iterative nature of the denoising process of the diffusion model. Future work could explore diffusion models optimized for faster inference to enable real-time expression code decoding while maintaining quality. Finally, our UHAP is trained on high-quality studio data with the same lighting conditions across subjects. Robustness to in-the-wild captures (e.g., mobile phone recordings under uncontrolled lighting) is still limited. Extending the UHAP training corpus to include light-stage data will increase robustness to such capture conditions.

Future Work. Future work will aim to address these limitations and further enhance our system’s capabilities. We plan to investigate methods for learning more natural and interactive gaze behaviors, perhaps by incorporating data from unscripted conversational videos or by enabling explicit gaze control. A significant avenue for future development involves leveraging our monocular expression encoder (E_{image}) by using it as an inference tool to collect expression codes on large-scale, diverse in-the-wild audio-visual datasets. This could enable the explicit modeling and audio-driven synthesis of a broader spectrum of nuanced human emotions, thereby enriching avatar expressiveness and realism.

REFERENCES

- Shivangi Aneja, Artem Sevastopolsky, Tobias Kirschstein, Justus Thies, Angela Dai, and Matthias Nießner. 2024a. GaussianSpeech: Audio-Driven Gaussian Avatars. arXiv:2411.18675 [cs.CV] <https://arxiv.org/abs/2411.18675>
- Shivangi Aneja, Justus Thies, Angela Dai, and Matthias Nießner. 2024b. FaceTalk: Audio-Driven Motion Diffusion for Neural Parametric Head Models. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*.
- Shivangi Aneja, Sebastian Weiss, Irene Baeza, Prashanth Chandran, Gaspard Zoss, Matthias Nießner, and Derek Bradley. 2025. ScaffoldAvatar: High-Fidelity Gaussian Avatars with Patch Expressions. arXiv:2507.10542 [cs.GR] <https://arxiv.org/abs/2507.10542>
- Monica Villanueva Aylagas, Hector Anadon Leon, Mattia Teye, and Konrad Tollmar. 2022. Voice2face: Audio-driven facial and tongue rig animations with cvaes.
- Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations. arXiv:2006.11477 [cs.CL] <https://arxiv.org/abs/2006.11477>
- 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 (SIGGRAPH ’99)*. 187–194. <https://doi.org/10.1145/311535.311556>
- Chen Cao, Tomas Simon, Jin Kyu Kim, Gabe Schwartz, Michael Zollhöfer, Shunsuke Saito, Stephen Lombardi, Shih-En Wei, Danielle Belko, Shouo-I Yu, Yaser Sheik, and Jason M. Saragih. 2022. Authentic volumetric avatars from a phone scan. *ACM Trans. Graph.* 41, 4 (2022), 163:1–163:19.
- Lele Chen, Zhiheng Li, Ross K. Maddox, Zhiyao Duan, and Chenliang Xu. 2018. Lip Movements Generation at a Glance. In *Computer Vision – ECCV 2018: 15th European Conference, Munich, Germany, September 8–14, 2018, Proceedings, Part VII* (Munich, Germany). Springer-Verlag, Berlin, Heidelberg, 538–553. https://doi.org/10.1007/978-3-030-01234-2_32
- J. S. Chung and A. Zisserman. 2016. Out of time: automated lip sync in the wild. In *Workshop on Multi-view Lip-reading, ACCV*.
- Radek Daněček, Michael J. Black, and Timo Bolkart. 2022. EMOCA: Emotion Driven Monocular Face Capture and Animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 20311–20322.
- Xiangyu Fan, Jiaqi Li, Zhiqian Lin, Weiyi Xiao, and Lei Yang. 2022a. FaceFormer: Speech-Driven 3D Facial Animation with Transformers. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 18770–18780. <https://doi.org/10.1109/CVPR52688.2022.01828>
- Yingruo Fan, Zhaojiang Lin, Jun Saito, Wenping Wang, and Taku Komura. 2022b. FaceFormer: Speech-Driven 3D Facial Animation with Transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

- 913 18749–18758. <https://doi.org/10.1109/CVPR52688.2022.01821>
- 914 Simon Giebenhain, Tobias Kirschstein, Martin Rünz, Lourdes Agapito, and Matthias
915 Nießner. 2024. NPGA: Neural Parametric Gaussian Avatars. In *SIGGRAPH Asia
916 Conference Papers (SA Conference Papers '24), December 3-6, Tokyo, Japan.*
<https://doi.org/10.1145/3680528.3687689>
- 917 Philip-William Grassal, Malte Prinzel, Titus Leistner, Carsten Rother, Matthias Nießner,
918 and Justus Thies. 2022. Neural Head Avatars from Monocular RGB Videos. In *IEEE
919 Conf. Comput. Vis. Pattern Recog.* IEEE, 18632–18643.
- 920 Jiazhui Guan, Zhanwang Zhang, Hang Zhou, Tianshu Hu, Kaisiyuan Wang, Dongliang
921 He, Haocheng Feng, Jingtuo Liu, Errui Ding, Ziwei Liu, and Jingdong Wang. 2023.
922 StyleSync: High-Fidelity Generalized and Personalized Lip Sync in Style-based
923 Generator. arXiv:2305.05445 [cs.CV] <https://arxiv.org/abs/2305.05445>
- 924 Jianzhu Guo, Dingyun Zhang, Xiaoqiang Liu, Zhizhou Zhong, Yuan Zhang, Pengfei
925 Wan, and Di Zhang. 2024. LivePortrait: Efficient Portrait Animation with Stitching
926 and Retargeting Control. *arXiv preprint arXiv:2407.03168* (2024).
- 927 Yudong Guo, Keyu Chen, Sen Liang, Yong-Jin Liu, Hujun Bao, and Juyong Zhang. 2021.
928 AD-NeRF: Audio Driven Neural Radiance Fields for Talking Head Synthesis. In
929 *IEEE/CVF International Conference on Computer Vision (ICCV)*. 5784–5793. <https://doi.org/10.1109/ICCV48922.2021.00572>
- 930 Yang Haotian, Zheng Mingwu, Ma ChongYang, Lai Yu-Kun, Wan Pengfei, and Huang
931 Haibin. 2024. VRMM: A Volumetric Relightable Morphable Head Model. In *SIG-
932 GRAPH 2024 Conference Proceedings*.
- 933 Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising Diffusion Probabilistic
934 Models. In *Advances in Neural Information Processing Systems*, H. Larochelle,
935 M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates,
936 Inc., 6840–6851. https://proceedings.neurips.cc/paper_files/paper/2020/file/4c5bcfec8584af0d967f1ab10179ca4b-Paper.pdf
- 937 Cheng hsin Wuu, Ningyuan Zheng, Scott Ardisson, Rohan Bali, Danielle Belko, Eric
938 Brockmeyer, Lucas Evans, Timothy Godisart, Hyowon Ha, Xuhua Huang, Alexander
939 Hypes, Taylor Koska, Steven Krenn, Stephen Lombardi, Xiaomin Luo, Kevyn
940 McPhail, Laura Millerschoen, Michal Perdoch, Mark Pitts, Alexander Richard, Jason
941 Saragih, Junko Saragih, Takaaki Shiratori, Tomas Simon, Matt Stewart, Autumn Trimble,
942 Xinshou Weng, David Whitewolf, Chenglei Wu, Shouo-I Yu, and Yaser Sheikh.
943 2023. Multiface: A Dataset for Neural Face Rendering. arXiv:2207.11243 [cs.CV]
944 <https://arxiv.org/abs/2207.11243>
- 945 Justin Johnson, Alexandre Alahi, and Li Fei-Fei. 2016. Perceptual Losses for Real-Time
946 Style Transfer and Super-Resolution. *CoRR* abs/1603.08155 (2016). arXiv:1603.08155
947 <http://arxiv.org/abs/1603.08155>
- 948 Tero Karras, Timo Aila, Samuli Laine, Antti Hervä, and Jaakko Lehtinen. 2017. Audio-
949 driven facial animation by joint end-to-end learning of pose and emotion.
- 950 Bernhard Kerbl, Georgios Kopanas, Thomas Leimkuhler, and George Drettakis. 2023b.
951 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Trans. Graph.*
952 42, 4, Article 139 (jul 2023), 14 pages. <https://doi.org/10.1145/3592433>
- 953 Thomas Kerbl, Luca Guarnera, Gerald Wimmer, Michael Wimmer, and Markus Stein-
954 berger. 2023a. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. In
955 *ACM SIGGRAPH Conference Proceedings*. <https://doi.org/10.1145/3588432.3591528>
- 956 Tobias Kirschstein, Javier Romero, Artem Sevastopolsky, Matthias Nießner, and Shunsuke
957 Saito. 2025. Avat3r: Large Animatable Gaussian Reconstruction Model for High-
958 fidelity 3D Head Avatars. arXiv:2502.20220 [cs.CV] <https://arxiv.org/abs/2502.20220>
- 959 Junxuan Li, Chen Cao, Gabriel Schwartz, Rawal Khirodkar, Christian Richardt, Tomas
960 Simon, Yaser Sheikh, and Shunsuke Saito. 2023. ER-NeRF: Efficient Region-Aware
961 Neural Radiance Fields for High-Fidelity Talking Portrait Synthesis. In *IEEE/CVF
962 International Conference on Computer Vision (ICCV)*.
- 963 Junxuan Li, Chen Cao, Gabriel Schwartz, Rawal Khirodkar, Christian Richardt, Tomas
964 Simon, Yaser Sheikh, and Shunsuke Saito. 2024. URAvatar: Universal Relightable
965 Gaussian Codec Avatars. In *ACM SIGGRAPH 2024 Conference Papers*.
- 966 Linzhou Li, Yumeng Li, Yanlin Weng, Youyi Zheng, and Kun Zhou. 2025a. RGBAv-
967 ator: Reduced Gaussian Blendshapes for Online Modeling of Head Avatars. In *The
968 IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- 969 Tianye Li, Timo Bolkart, Michael J. Black, Hao Li, and Javier Romero. 2017. Learning a
970 model of facial shape and expression from 4D scans. *ACM Trans. Graph.* 36, 6 (2017),
971 194:1–194:17.
- 972 Xuanchen Li, Jianyu Wang, Yuhao Cheng, Yikun Zeng, Xingyu Ren, Wenhan Zhu,
973 Weiming Zhao, and Yichao Yan. 2025b. Towards High-fidelity 3D Talking Avatar
974 with Personalized Dynamic Texture. *arXiv preprint arXiv:2503.00495* (2025).
- 975 Shengjie Ma, Yanlin Weng, Tianjia Shao, and Kun Zhou. 2024. 3D Gaussian Blendshapes
976 for Head Avatar Animation. In *ACM SIGGRAPH Conference Proceedings, Denver, CO,
977 United States, July 28 - August 1, 2024*.
- 978 Julieta Martinez, Emily Kim, Javier Romero, Timur Bagautdinov, Shunsuke Saito, Shouo-
979 i Yu, Stuart Anderson, Michael Zollhöfer, Te-Li Wang, Shaojin Bai, Chenghui Li, Shih-
980 En Wei, Rohan Joshi, Wyatt Borsos, Tomas Simon, Jason Saragih, Paul Theodosis,
981 Alexander Greene, Anjani Josyula, Silvio Mano Maeta, Andrew I. Jewett, Simon
982 Venshtain, Christopher Heilman, Yueh-Tung Chen, Sidi Fu, Mohamed Ezzeldin A.
983 Elshaer, Tingfang Du, Longhua Wu, Shen-Chi Chen, Kai Kang, Michael Wu, Youssef
984 Emad, Steven Longay, Ashley Brewer, Hitesh Shah, James Booth, Taylor Koska,
985 Kayla Haidle, Matt Andromalos, Joanna Hsu, Thomas Dauer, Peter Selednik, Tim
986 Godisart, Scott Ardisson, Matthew Cipperly, Ben Humberston, Lon Farr, Bob Hansen,
987 Peihong Guo, Dave Braun, Steven Krenn, He Wen, Lucas Evans, Natalia Fadeeva,
988 Matthew Stewart, Gabriel Schwartz, Divam Gupta, Gyeongsik Moon, Kaiwen Guo,
989 Yuan Dong, Yichen Xu, Takaaki Shiratori, Fabian Prada, Bernardo R. Pires, Bo Peng,
990 Julia Buffalini, Autumn Trimble, Kevyn McPhail, Melissa Schoeller, and Yaser Sheikh.
991 2024. Codec Avatar Studio: Paired Human Captures for Complete, Driveable, and
992 Generalizable Avatars. *NeurIPS Track on Datasets and Benchmarks* (2024).
- 993 Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ra-
994 mamoorthi, and Ren Ng. 2020. NeRF: Representing Scenes as Neural Radiance Fields
995 for View Synthesis. In *Proceedings of the European Conference on Computer Vision
996 (ECCV)*. 405–421. https://doi.org/10.1007/978-3-030-58452-8_24
- 997 Evonne Ng, Javier Romero, Timur Bagautdinov, Shaojie Bai, Trevor Darrell, Angjoo
998 Kanazawa, and Alexander Richard. 2024. From Audio to Photoreal Embodiment:
999 Synthesizing Humans in Conversations. arXiv:2401.01885 [cs.CV] <https://arxiv.org/abs/2401.01885>
- 999 Ziqiao Peng, Haoyu Wu, Zhenbo Song, Hao Xu, Xiangyu Zhu, Jun He, Hongyan Liu,
999 and Zhaoxin Fan. 2023. Emotalk: Speech-driven emotional disentanglement for 3d
999 face animation.
- 999 Ethan Perez, Florian Strub, Harm de Vries, Vincent Dumoulin, and Aaron C. Courville.
999 2018. FiLM: Visual Reasoning with a General Conditioning Layer. In *AAAI*.
999 Shenhan Qian, Tobias Kirschstein, Liam Schoneveld, Davide Davoli, Simon Giebenhain,
999 and Matthias Nießner. 2024a. Gaussianavatars: Photorealistic head avatars with
999 rigged 3d gaussians. In *Proceedings of the IEEE/CVF Conference on Computer Vision
999 and Pattern Recognition*. 20299–20309.
- 999 Shenhan Qian, Tobias Kirschstein, Liam Schoneveld, Davide Davoli, Simon Giebenhain,
999 and Matthias Nießner. 2024b. GaussianAvatars: Photorealistic Head Avatars with
999 Rigged 3D Gaussians. In *Proceedings of the IEEE/CVF Conference on Computer Vision
999 and Pattern Recognition (CVPR)*. https://openaccess.thecvf.com/content/CVPR2024/html/Qian_GaussianAvatars_Photorealistic_Head_Avatars_with_Rigged_3D_Gaussians_CVPR_2024_paper.html
- 999 Alexander Richard, Colin Lea, Shugao Ma, Jurgen Gall, Fernando de la Torre, and Yaser
999 Sheikh. 2021a. Audio- and Gaze-Driven Facial Animation of Codec Avatars. In
999 *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision
999 (WACV)*. 41–50.
- 999 Alexander Richard, Michael Zollhöfer, Yandong Wen, Fernando De la Torre, and Yaser
999 Sheikh. 2021b. MeshTalk: 3D Face Animation from Speech Using Cross-Modality
999 Disentanglement. In *IEEE/CVF International Conference on Computer Vision (ICCV)*.
999 1173–1182. <https://doi.org/10.1109/ICCV48922.2021.00121>
- 999 Shunsuke Saito, Gabriel Schwartz, Tomas Simon, Junxuan Li, and Giljoo Nam. 2024.
999 Relightable Gaussian Codec Avatars. In *CVPR*.
- 999 Stefan Stan, Kazi Injamamul Haque, and Zerrin Yumak. 2023a. FaceDiffuser: Speech-
999 Driven 3D Facial Animation Synthesis Using Diffusion. arXiv:2309.11306 [cs.CV]
999 <https://arxiv.org/abs/2309.11306>
- 999 Stefan Stan, Kazi Injamamul Haque, and Zerrin Yumak. 2023b. FaceDiffuser: Speech-
999 Driven 3D Facial Animation Synthesis Using Diffusion. arXiv:2309.11306 [cs.CV]
999 <https://arxiv.org/abs/2309.11306>
- 999 Zhiyao Sun, Tian Lv, Sheng Ye, Matthieu Lin, Jenny Sheng, Yu-Hui Wen, Minjing Yu,
999 and Yong-Jin Liu. 2024a. DiffPoseTalk: Speech-Driven Stylistic 3D Facial Animation
999 and Head Pose Generation via Diffusion Models. *ACM Transactions on Graphics*
999 (2024). <https://doi.org/10.1145/3679561> Proceedings of SIGGRAPH 2024.
- 999 Zhiyao Sun, Tian Lv, Sheng Ye, Matthieu Lin, Jenny Sheng, Yu-Hui Wen, Minjing Yu,
999 and Yong-Jin Liu. 2024b. DiffPoseTalk: Speech-Driven Stylistic 3D Facial Animation
999 and Head Pose Generation via Diffusion Models. *ACM Transactions on Graphics
999 (TOG)* 43, 4, Article 46 (2024), 9 pages. <https://doi.org/10.1145/3658221>
- 999 Jiaxiang Tang, Kaisiyuan Wang, Hang Zhou, Xiaokang Chen, Dongliang He, Jingtuo
999 Liu, Tianshu Hu, Gang Zeng, and Jingdong Wang. 2022. RAD-NeRF: Real-Time
999 Neural Radiance Talking Portrait Synthesis via Audio-Spatial Decomposition. *arXiv
999 preprint arXiv:2211.12368* (2022).
- 999 Sarah Taylor, Taehwan Kim, Yisong Yue, Moshe Mahler, James Krahe, Anastasio Garcia
999 Rodriguez, Jessica Hodgins, and Iain Matthews. 2017. A Deep Learning Approach
999 for Generalized Speech Animation. *ACM Transactions on Graphics* 36, 4 (2017),
999 93:1–93:12. <https://doi.org/10.1145/3072959.3073699>
- 999 Kartik Teotia, Hyeongwoo Kim, Pablo Garrido, Marc Habermann, Mohamed Elgharib,
999 and Christian Theobalt. 2024. GaussianHeads: End-to-End Learning of Drivable
999 Gaussian Head Avatars from Coarse-to-fine Representations. *ACM Trans. Graph.*
999 43, 6, Article 264 (Nov. 2024), 12 pages. <https://doi.org/10.1145/3687927>
- 999 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N.
999 Gomez, Łukasz Kaiser, and Illia Polosukhin. 2023. Attention Is All You Need.
999 arXiv:1706.03762 [cs.CL] <https://arxiv.org/abs/1706.03762>
- 999 Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. 2004. Image
999 quality assessment: from error visibility to structural similarity. *IEEE Trans. Image
999 Process.* 13, 4 (2004), 600–612.
- 999 Sebastian Winberg, Gaspard Zoss, Prashanth Chandran, Paulo Gotardo, and Derek
999 Bradley. 2022. Facial hair tracking for high fidelity performance capture. *ACM Trans.
999 Graph.* 41, 4, Article 165 (July 2022), 12 pages. <https://doi.org/10.1145/3528223.3530116>

1027	Jinbo Xing, Menghan Xia, Yuechen Zhang, Xiaodong Cun, Jue Wang, and Tien-Tsin Wong. 2023. Codetalker: Speech-driven 3d facial animation with discrete motion prior.	1084
1028		1085
1029	Sicheng Xu, Guojun Chen, Yu-Xiao Guo, Jiaolong Yang, Chong Li, Zhenyu Zang, Yizhong Zhang, Xin Tong, and Baining Guo. 2024a. VASA-1: Lifelike Audio-Driven Talking Faces Generated in Real Time. arXiv:2404.10667 [cs.CV] https://arxiv.org/abs/2404.10667	1086
1030		1087
1031	Yuelang Xu, Lizhen Wang, Zerong Zheng, Zhaoqi Su, and Yebin Liu. 2024b. 3D Gaussian Parametric Head Model. In <i>Proceedings of the European Conference on Computer Vision (ECCV)</i> .	1088
1032		1089
1033	Zhenhui Ye, Ziyu Jiang, Yi Ren, Jinglin Liu, Jinzheng He, and Zhou Zhao. 2023. GeneFace: Generalized and High-Fidelity Audio-Driven 3D Talking Face Synthesis. <i>International Conference on Learning Representations (ICLR)</i> (2023).	1090
1034		1091
1035	Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In <i>CVPR</i> .	1092
1036		1093
1037		1094
1038		1095
1039		1096
1040		1097
1041		1098
1042		1099
1043		1100
1044		1101
1045		1102
1046		1103
1047		1104
1048		1105
1049		1106
1050		1107
1051		1108
1052		1109
1053		1110
1054		1111
1055		1112
1056		1113
1057		1114
1058		1115
1059		1116
1060		1117
1061		1118
1062		1119
1063		1120
1064		1121
1065		1122
1066		1123
1067		1124
1068		1125
1069		1126
1070		1127
1071		1128
1072		1129
1073		1130
1074		1131
1075		1132
1076		1133
1077		1134
1078		1135
1079		1136
1080		1137
1081		1138
1082		1139
1083		1140

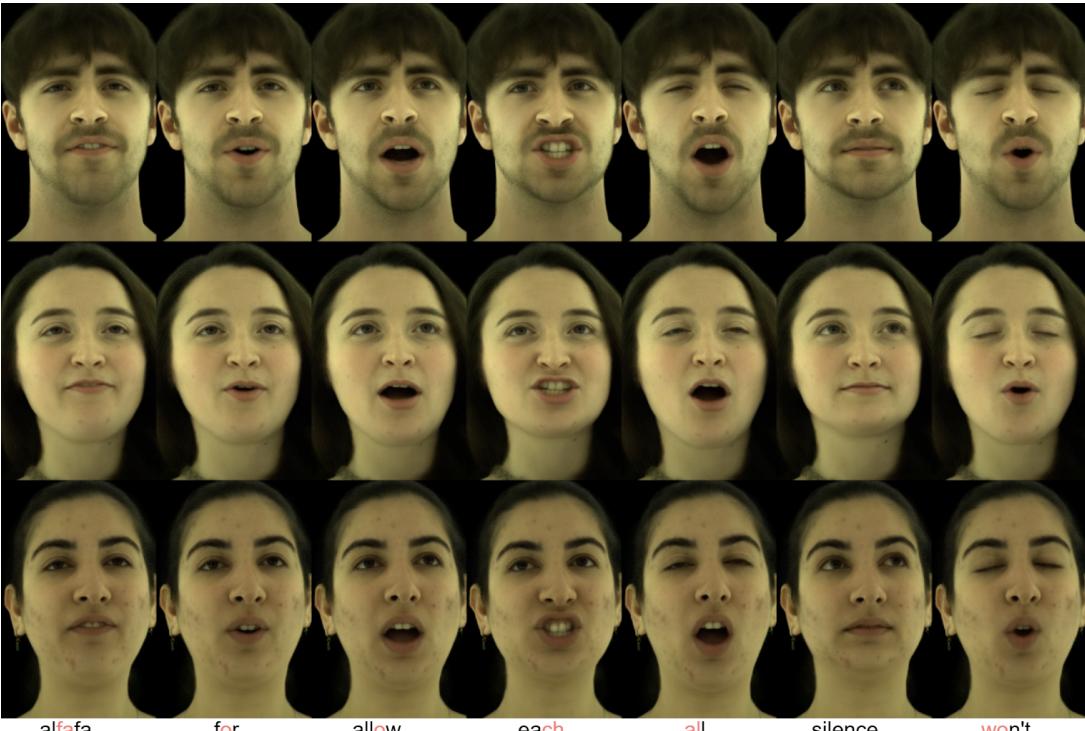


Fig. 4. Audio-driven synthesis results for three UHAP model identities with corresponding audio prompts.

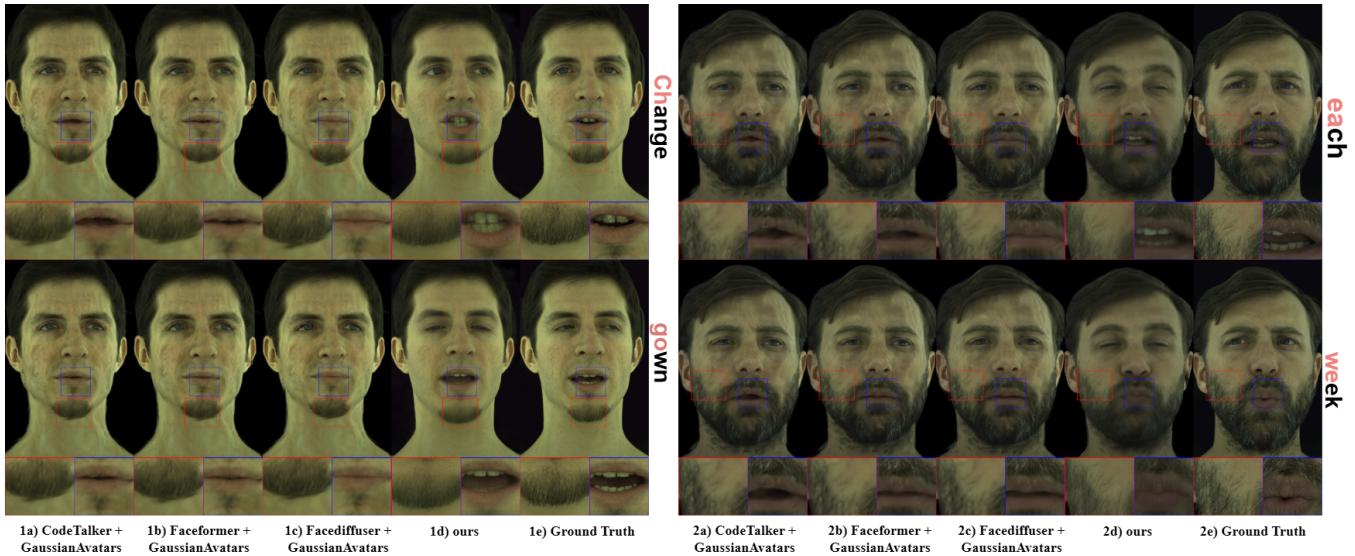


Fig. 5. Qualitative comparison with SOTA methods (CodeTalker+GA, Faceformer+GA, FaceDiffuser+GA, Ours) and Ground Truth for specified audio segments. GA denotes GaussianAvatars augmentation.

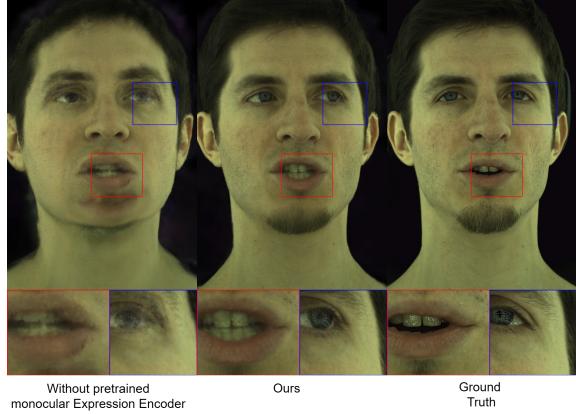


Fig. 8. Ablation on monocular encoder (E_{image}) training: Encoder trained during fitting (left), Ours (pretrained encoder, center), Ground Truth (right).

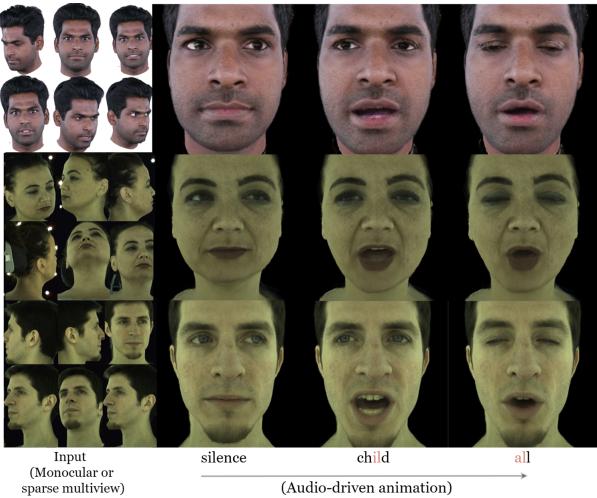


Fig. 9. Personalization from sparse inputs (left column) and resulting audio-driven animation (right columns) for three subjects at a novel viewpoint.

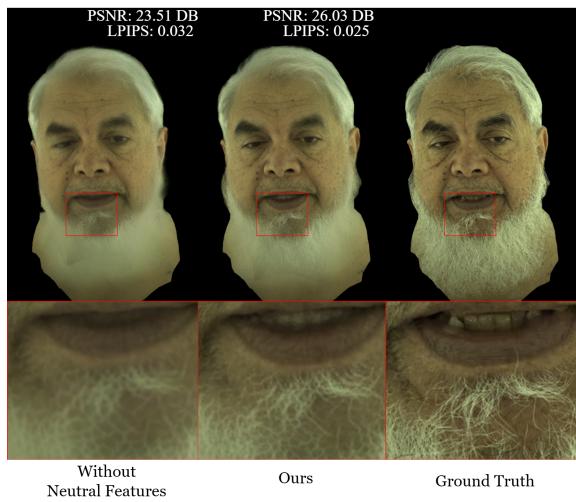


Fig. 10. Ablation on neutral features (f_{neut}): Without Neutral Features (left), Ours (center), Ground Truth (right), with PSNR/LPIPS metrics.

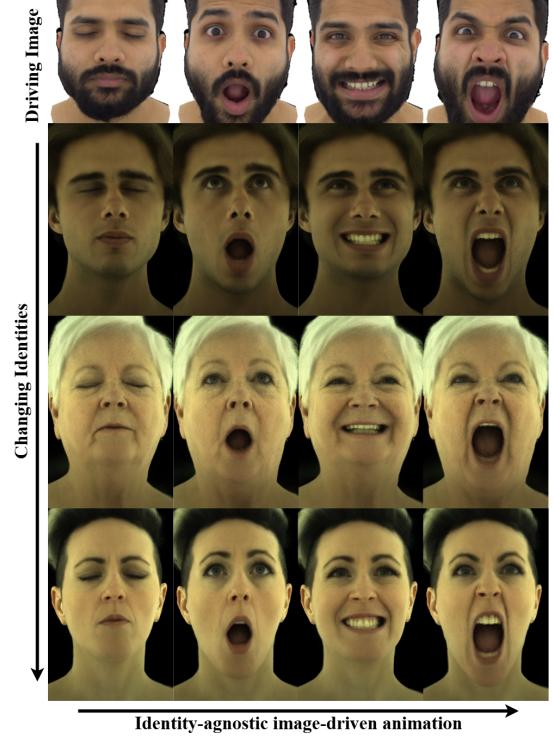


Fig. 11. Image-driven animation: Driving image sequence (top) and animated target identities (rows below).

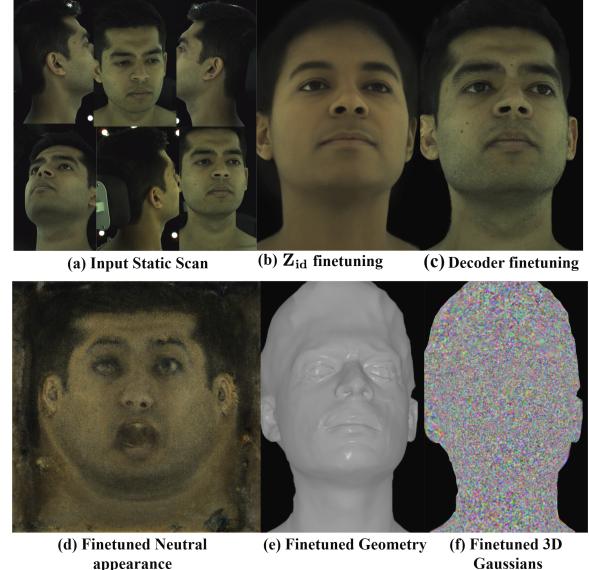


Fig. 12. Personalization pipeline stages: (a) Input static scan. (b) Result after Z_{id} fine-tuning. (c) Result after \mathcal{D}_{UHAP} decoder fine-tuning. Process yields (d) finetuned neutral appearance, (e) geometry, and (f) 3D Gaussians.