

GazeGaussian: High-Fidelity Gaze Redirection with 3D Gaussian Splatting

Anonymous CVPR submission

Paper ID 3250

Abstract

001 Gaze estimation encounters generalization challenges when
 002 dealing with out-of-distribution data. To address this prob-
 003 lem, recent methods use neural radiance fields (NeRF) to
 004 generate augmented data. However, existing methods based
 005 on NeRF are computationally expensive and lack facial de-
 006 tails. 3D Gaussian Splatting (3DGS) has become the prevail-
 007 ing representation of neural fields. While 3DGS has been
 008 extensively examined in head avatars, it faces challenges
 009 with accurate gaze control and generalization across dif-
 010 ferent subjects. In this work, we propose GazeGaussian, a
 011 high-fidelity gaze redirection method that uses a two-stream
 012 3DGS model to represent the face and eye regions separately.
 013 By leveraging the unstructured nature of 3DGS, we develop
 014 a novel eye representation for rigid eye rotation based on
 015 the target gaze direction. To enhance synthesis generaliza-
 016 tion across various subjects, we integrate an expression-
 017 conditional module to guide the neural renderer. Compre-
 018 hensive experiments show that GazeGaussian outperforms
 019 existing methods in rendering speed, gaze redirection ac-
 020 curacy, and facial synthesis across multiple datasets. We
 021 also demonstrate that existing gaze estimation methods can
 022 leverage GazeGaussian to improve their generalization per-
 023 formance. The code will be released.

024 1. Introduction

025 Gaze estimation is a fundamental component across various
 026 applications [1, 25, 27], yet current estimators [3, 4, 48] often
 027 struggle to generalize effectively to out-of-distribution data.
 028 To address this, recent approaches [34, 40, 51] have started
 029 exploring gaze redirection, which manipulates the gaze in an
 030 input image toward a target direction. This process generates
 031 augmented data to enhance the generalization capabilities of
 032 gaze estimators.

033 Earlier methods [10, 52, 53, 56] formulate gaze redirection
 034 as a 2D image manipulation task, relying on deep learn-
 035 ing techniques to warp eye regions of the image toward
 036 the target gaze direction. However, these 2D approaches
 037 overlook the inherently 3D nature of head and gaze ma-

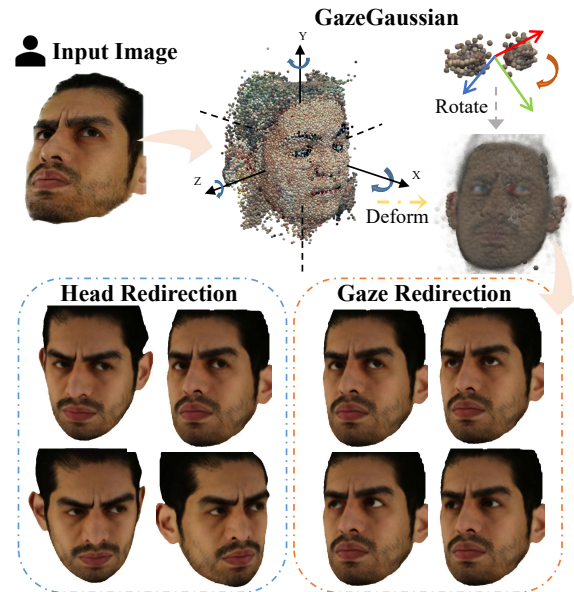


Figure 1. GazeGaussian for gaze redirection: Given an input image, GazeGaussian deforms face and eye Gaussians from canonical space to generate high-fidelity head images with accurate gaze redirection.

nipulation, often resulting in poor spatial consistency and
 limited synthesis fidelity. With advancements in Neural
 Radiance Fields (NeRF) [26] and its variants [42, 44], sev-
 eral methods [12, 16, 59, 61] have achieved 3D dynamic
 head representation and high-fidelity avatar synthesis. Mean-
 while, to enable precise control of gaze direction, recent
 research [34, 40, 51] has introduced approaches that decou-
 ple the face and eye regions, modeling each with separate
 neural fields to achieve accurate gaze redirection.

As NeRF-based methods are hindered by high compu-
 tational demands, 3D Gaussian Splatting [18] and its
 variants [17, 24, 43] achieve impressive rendering qual-
 ity with significantly faster training speeds. Recent re-
 search [31, 47, 50] has applied these methods to 3D head
 animation, typically using face-tracking [39, 60] param-
 eters to model dynamic 3D head representations. However,
 existing 3DGS-based approaches neglect the accurate control

of gaze direction and struggle to generalize across different subjects, limiting their effectiveness for gaze redirection.

To address the above issues, we propose GazeGaussian, a high-fidelity gaze redirection method that leverages a two-stream 3D Gaussian Splatting (3DGS) model to represent the face and eye regions, respectively. To the best of our knowledge, this is the first integration of 3DGS into gaze redirection tasks. An overview is shown in Fig. 1.

GazeGaussian begins by initializing the two-stream 3DGS model using a pre-trained neutral mesh on the training dataset. This mesh is divided into distinct regions for the face and eyes. By employing gaze direction and face tracking codes, we optimize a deformation field for the face and a rotation field for the eyes, allowing us to adjust the neutral Gaussians accordingly. To achieve precise eye rotation aligned with the target gaze, we present a novel Gaussian Eye Rotation Representation (GERR). In contrast to methods like GazeNeRF that implicitly alter feature maps, GazeGaussian explicitly adjusts the position of Gaussians in the eye branch according to the desired gaze direction, utilizing the controllable nature of 3DGS. To address possible errors in gaze direction, GazeGaussian develops an eye rotation field to enhance redirection accuracy. The two-stream Gaussians are rasterized into high-level features and sent to the neural renderer. Finally, to enhance synthesis generalization across different subjects and preserve facial details, we employ an expression-guided neural renderer (EGNR) to synthesize the final gaze-redirection images.

Our main contributions are summarized as follows:

- We introduce GazeGaussian, the first 3DGS-based gaze redirection pipeline, achieving precise gaze manipulation and high-fidelity head avatar synthesis.
- To enable rigid and accurate eye rotation based on the target gaze direction, we propose a novel two-stream 3DGS framework to decouple face and eye deformations, featuring a specialized Gaussian eye rotation for explicit control over eye movement.
- To enhance the synthesis generalization of 3DGS, we design an expression-guided neural renderer (EGNR) to retain facial details across various subjects.
- We conduct comprehensive experiments on ETH-XGaze, ColumbiaGaze, MPIIFaceGaze, and GazeCapture datasets, where GazeGaussian achieves state-of-the-art gaze redirection accuracy and facial synthesis quality with competitive rendering speed.

2. Related Work

Gaze Redirection. Gaze redirection is the task of manipulating the gaze direction of a face image to a target direction while preserving the subject’s identity and other facial details. Earlier approaches for gaze redirection include novel view synthesis [5, 11, 21], eye-replacement [32, 36], and warping-based methods [10, 19, 45]. However, these methods are

limited by person-specific data requirements, restricted redirection range, and artifact introduction. To further improve gaze redirection, recent studies [14, 28, 46, 58] have employed neural network-based generative models. STED [58], building on the FAZE [28], introduces a self-transforming encoder-decoder that generates full-face images with high-fidelity control over gaze direction and head pose. Effective gaze redirection should account for both the 3D nature of eyeball rotation and the deformation of surrounding facial regions. With advancements in Neural Radiance Fields (NeRF) [26], several studies [22, 34, 40, 51] have aimed to model the complex rotation of the eyeball. GazeNeRF [34] employs a two-stream MLP architecture to separately model the face only and eye regions, achieving improved gaze redirection performance.

However, these methods are hindered by substantial computational demands and limited rendering efficiency. Additionally, gaze manipulation occurs at the feature map level and remains an implicit approach. In contrast, GazeGaussian allows for explicit control over eye rotations, improving gaze redirection accuracy and accelerating the synthesis process.

Head Avatar Synthesis. The synthesis of head avatars has garnered considerable attention in recent years. FLAME [23] is a parameterized 3D head model that maps parameters of shape, expression, and pose onto 3D facial geometry, allowing for realistic and controllable head avatar generation. Many subsequent works [2, 6, 8, 29, 30, 33] focus on using the FLAME model for speech-driven head avatar animation. Recent head animation techniques can be categorized into two main approaches: NeRF-based methods and 3DGS-based methods. NeRF-based approaches [9, 16, 59, 61] leverage neural radiance fields to deform facial movements from a canonical space. HeadNeRF [16] introduces a parametric head model that controls facial shape, expression, and albedo under different lighting conditions. With the emergence of 3D Gaussian Splatting (3DGS) [18], several approaches [7, 31, 47, 50] have explored its application in head avatar modeling. Gaussian Head Avatar [49] initializes Gaussians with a neutral mesh head and incorporates MLPs to deform complex facial expressions.

While these methods produce impressive results in creating 3D head avatars, they overlook precise gaze control and do not generalize well across different subjects. In contrast, GazeGaussian emphasizes precise gaze direction control by decoupling facial animations and gaze movement within a two-stream model. Furthermore, we introduce an expression-guided neural renderer designed to improve the quality of synthesis.

3. Overview

The pipeline of GazeGaussian is illustrated in Fig. 2, including the two-stream Gaussians and the proposed expression-guided neural renderer. Before the beginning of the pipeline,

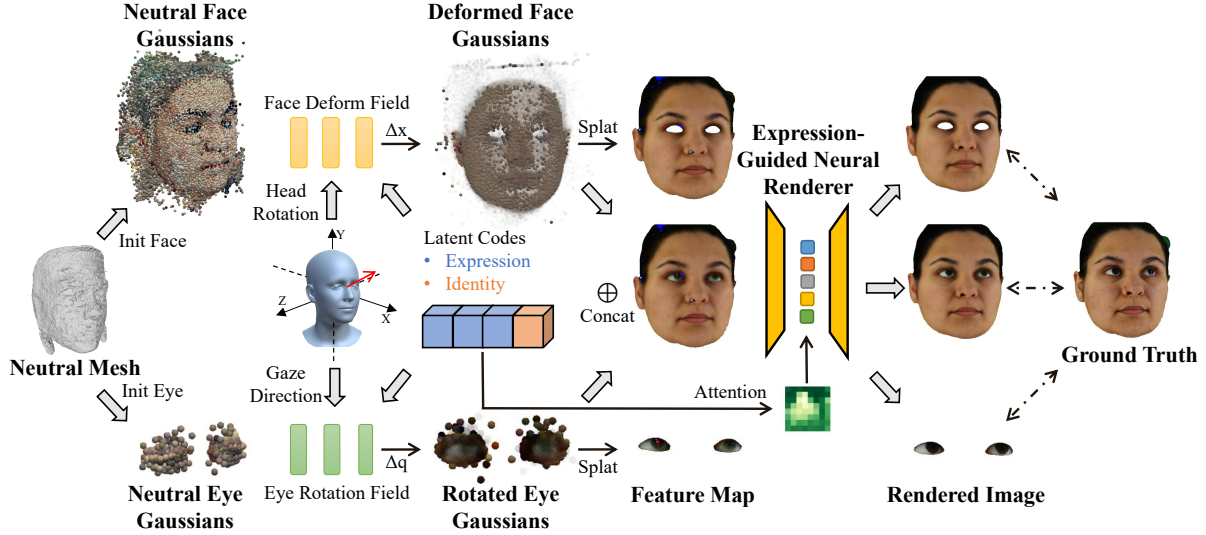


Figure 2. Pipeline of GazeGaussian. We initialize face-only and eye regions from a pre-trained neutral mesh. Using target expression codes, head rotation, and gaze direction, GazeGaussian optimizes face deformation and eye rotation fields to transform the neutral Gaussians. The transformed Gaussians are splatted into feature maps. The expression codes guide the neural renderer through cross-attention, enabling the rendering of feature maps into high-fidelity images, which are then supervised by multi-view RGB images.

we follow the data preprocessing in GazeNeRF [34] and Gaussian Head Avatar [50], which include background removal, gaze direction normalization, and facial tracking for each frame. To obtain a neutral mesh for Gaussian initialization, we first reconstruct a Sign Distance Function (SDF) based neutral geometry and then optimize a face deformation field and an eye rotation field from the training data. A neutral mesh representing a coarse geometry across different subjects can be extracted using DM Tet [35]. We then partition the neutral mesh into face-only and eye regions using 3D landmarks, initializing the two-stream Gaussians. Based on these neutral Gaussians, GazeGaussian optimizes a face deformation field and an eye rotation field to transform the Gaussians according to the target expression codes, gaze direction, and head rotation. Next, we concatenate the two-stream Gaussians and rasterize them into a high-dimensional feature map representing the head, face-only, and eye regions. Finally, these feature maps are fed into the expression-guided neural renderer to generate high-fidelity gaze redirection images. The ground truth image is used to supervise the rendered face-only, head, and eye images.

4. Method

4.1. Preliminaries

The vanilla 3D Gaussians [18] with N points are represented by their positions X , the multi-channel color C , the rotation Q , scale S and opacity A . The color C is computed using spherical harmonics, and the rotation Q is represented as the quaternion. These Gaussians are then rasterized and rendered to a multi-channel image I based on the camera

parameters μ . This rendering process can be expressed as:

$$I = \mathcal{R}(X, C, Q, S, A; \mu), \quad (1)$$

4.2. Two-stream GazeGaussian Representation

Our task is to synthesize a head avatar conditioned on gaze direction, head rotation, and expression latent codes. To decouple the complex movements in the face and eyes, we introduce a two-stream Gaussian model consisting of a face-only branch and an eye branch. In the following subsections, we will describe the face deformation and eye rotation processes, respectively.

4.2.1. Face Deformation

For the face-only branch, inspired by Gaussian Head Avatar, we first construct canonical neutral face Gaussians with attributes: $\{X_0^f, F_0^f, Q_0^f, S_0^f, A_0^f\}$, which are fully optimizable. $X_0^f \in \mathbb{R}^{N \times 3}$ represents the positions of the Gaussians in the canonical space. $F_0^f \in \mathbb{R}^{N \times 128}$ denotes the point-wise feature vectors as their intrinsic properties. $Q_0^f \in \mathbb{R}^{N \times 4}$, $S_0^f \in \mathbb{R}^{N \times 3}$ and $A_0^f \in \mathbb{R}^{N \times 1}$ denotes the neutral rotation, scale and opacity respectively. The neutral color is directly predicted from the point-wise feature vectors F_0^f . Then we construct several MLPs, denoted as Φ^f , to serve as face deformation fields that transform the neutral face Gaussians. Next, we describe the process of applying offsets to each Gaussian attribute.

Positions X^f of the Gaussians. We predict the displacements respectively controlled by the latent codes and the head pose in the canonical space through two different MLPs:

215 ${}_{def}^{exp} \mathcal{F}^f \in \Phi^f$ and ${}_{def}^{pose} \mathcal{F}^f \in \Phi^f$. Then, we add them to the
216 neutral positions.

$$217 \quad X^f = \mathbf{X}_0^f + \lambda_{exp}(\mathbf{X}_0^f) {}_{def}^{exp} \mathcal{F}^f(\mathbf{X}_0^f, \theta) \quad (2)$$

$$218 \quad + \lambda_{pose}(\mathbf{X}_0^f) {}_{def}^{pose} \mathcal{F}^f(\mathbf{X}_0^f, \beta),$$

218 θ denoting latent codes including expression and identity co-
219 efficients and β denoting the head pose. $\lambda_{exp}(\cdot)$ and $\lambda_{pose}(\cdot)$
220 represent the degree to which the point is influenced by the
221 expression or head pose, respectively, which can be calcu-
222 lated as:

$$223 \quad \lambda_{exp}(x) = \begin{cases} 1, & dist(x, \mathbf{P}_0^f) < t_1 \\ \frac{t_2 - dist(x, \mathbf{P}_0^f)}{t_2 - t_1}, & dist(x, \mathbf{P}_0^f) \in [t_1, t_2] \\ 0, & dist(x, \mathbf{P}_0^f) > t_2 \end{cases}$$

224 with $\lambda_{pose}(x) = 1 - \lambda_{exp}(x)$, where $x \in \mathbf{X}_0^f$ denotes the
225 position of a neutral Gaussian, $dist(x, \mathbf{P}_0^f)$ represents the
226 minimum distance from point x to the 3D landmarks (with-
227 out eyes) \mathbf{P}_0^f . Following the approach in Gaussian Head
228 Avatar, the predefined hyperparameters are set as $t_1 = 0.15$
229 and $t_2 = 0.25$.

230 **Color** C^f of the Gaussians. Modeling dynamic details
231 requires a color that varies with expressions. The color
232 is directly predict by two color MLPs: ${}_{col}^{exp} \mathcal{F}^f \in \Phi^f$ and
233 ${}_{col}^{pose} \mathcal{F}^f \in \Phi^f$:

$$234 \quad C^f = \lambda_{exp}(\mathbf{X}_0^f) {}_{col}^{exp} \mathcal{F}^f(\mathbf{F}_0^f, \theta) \quad (3)$$

$$235 \quad + \lambda_{pose}(\mathbf{X}_0^f) {}_{col}^{pose} \mathcal{F}^f(\mathbf{F}_0^f, \beta),$$

235 **Rotation, Scale and Opacity** $\{Q^f, S^f, A^f\}$ of the Gaus-
236 sians. These three attributes are also dynamic, capturing
237 detailed expression-related appearance changes. We just use
238 another two attribute MLPs: ${}_{att}^{exp} \mathcal{F}^f \in \Phi^f$ and ${}_{att}^{pose} \mathcal{F}^f \in \Phi^f$
239 to predict their shift from the neutral value.

$$240 \quad \{Q^f, S^f, A^f\} = \{\mathbf{Q}_0^f, \mathbf{S}_0^f, \mathbf{A}_0^f\} \quad (4)$$

$$241 \quad + \lambda_{exp}(\mathbf{X}_0^f) {}_{att}^{exp} \mathcal{F}^f(\mathbf{F}_0^f, \theta)$$

$$242 \quad + \lambda_{pose}(\mathbf{X}_0^f) {}_{att}^{pose} \mathcal{F}^f(\mathbf{F}_0^f, \beta),$$

241 Finally, we apply rigid rotations and translations to trans-
242 form Gaussians in the canonical space to the world space.
243 Then, these Gaussians are rasterized into the feature maps.
244 The above face-only branch can be formulated as:

$$245 \quad \mathcal{M}_f = \mathcal{R}(\{X^f, C^f, Q^f, S^f, A^f\}) \quad (5)$$

$$246 \quad = \mathcal{R}(\Phi^f(\mathbf{X}_0^f, \mathbf{F}_0^f, \mathbf{Q}_0^f, \mathbf{S}_0^f, \mathbf{A}_0^f; \theta, \beta)),$$

246 where \mathcal{R} represents the rasterizer and \mathcal{M}_f indicates the
247 feature map from the face-only branch.

4.2.2. Eye Rotation

248 For the eye branch, we also construct canonical neutral eye
249 Gaussians with attributes $\{\mathbf{X}_0^e, \mathbf{F}_0^e, \mathbf{Q}_0^e, \mathbf{S}_0^e, \mathbf{A}_0^e\}$. These
250 attributes share the same dimensionality as those in the face-
251 only branch, except that $\mathbf{S}_0^e \in \mathbb{R}^{N \times 1}$ is constrained to be
252 spherical, aligning with the rotational properties of the eye-
253 ball. Next, we describe the process of applying offsets to
254 each Gaussian attribute.

255 **Positions** X^e of the Gaussians. Directly applying the
256 same deformation strategy as for the face branch would
257 fail to fully leverage the unique characteristics of eyeball
258 rotational motion, resulting in insufficient gaze redirection
259 accuracy. Therefore, we first rotate the eye Gaussians in
260 the canonical space and then incorporate the eye geometry
261 information from the latent codes of different subjects to
262 generate biases. Since the gaze labels may contain noise,
263 directly using the normalized gaze direction φ to rotate the
264 Gaussians would lead to numerical optimization errors. To
265 address this, we optimize two separate MLPs: ${}_{rot}^{gaze} \mathcal{F}^e \in \Phi^e$
266 and ${}_{def}^{exp} \mathcal{F}^e \in \Phi^e$ to predict the biases for Gaussian rotation
267 and displacement.

$$268 \quad X^e = {}_{def}^{exp} \mathcal{F}^e(\mathbf{X}_0^e, \theta) + {}_{rot}^{gaze} \mathcal{F}^e(\mathbf{X}_0^e, \varphi) \mathbf{X}_0^e, \quad (6)$$

270 Since eyes are relatively small and mainly influenced by the
271 gaze direction, λ used in the face is omitted here.

272 **Color** C^e of the Gaussians. The color of the eye region
273 is influenced by the gaze direction and latent codes. We use
274 two MLPs: ${}_{col}^{exp} \mathcal{F}^e \in \Phi^e$ and ${}_{col}^{gaze} \mathcal{F}^e \in \Phi^e$ to predict it:

$$275 \quad C^e = {}_{att}^{exp} \mathcal{F}^e(\mathbf{F}_0^e, \theta) + {}_{col}^{gaze} \mathcal{F}^e(\mathbf{X}_0^e, \varphi), \quad (7)$$

276 **Rotation, Scale and Opacity** $\{Q^e, S^e, A^e\}$ of the Gaus-
277 sians. We just use another two attribute MLPs ${}_{att}^{exp} \mathcal{F}^e \in \Phi^e$
278 and ${}_{att}^{gaze} \mathcal{F}^e \in \Phi^e$ to predict their shift.

$$279 \quad \{Q^e, S^e, A^e\} = \{\mathbf{Q}_0^e, \mathbf{S}_0^e, \mathbf{A}_0^e\} + {}_{att}^{exp} \mathcal{F}^e(\mathbf{F}_0^e, \theta) \quad (8)$$

$$280 \quad + {}_{att}^{gaze} \mathcal{F}^e(\mathbf{F}_0^e, \varphi),$$

280 Finally, we transform Gaussians in the canonical space
281 to the world space. Then these eye Gaussians are rasterized
282 into the feature maps. The eye branch is formulated as:

$$283 \quad \mathcal{M}_e = \mathcal{R}(\{X^e, C^e, Q^e, S^e, A^e\}) \quad (9)$$

$$284 \quad = \mathcal{R}(\Phi^e(\mathbf{X}_0^e, \mathbf{F}_0^e, \mathbf{Q}_0^e, \mathbf{S}_0^e, \mathbf{A}_0^e; \theta, \varphi)),$$

284 To obtain the full head rendering, we simply concat the
285 two-stream Gaussians and rasterized them into feature maps:

$$286 \quad \mathcal{M}_h = \mathcal{R}(\{X^f, C^f, Q^f, S^f, A^f\} \quad (10)$$

$$287 \quad \{X^e, C^e, Q^e, S^e, A^e\}),$$

4.3. Expression-Guided Neural Renderer

288 After obtaining the rasterized feature maps from Gaussians,
289 a UNet-like neural renderer \mathbf{R} opts to synthesize the final
290

291 face-only, eyes, and head images $\{\mathcal{I}_f, \mathcal{I}_e, \mathcal{I}_h\}$:

$$292 \quad \{\mathcal{I}_f, \mathcal{I}_e, \mathcal{I}_h\} = \mathbf{R}(\{\mathcal{M}_f, \mathcal{M}_e, \mathcal{M}_h\}), \quad (11)$$

293 To enhance the generalization ability across different sub-
294 jects, we inject the latent codes θ into the neural renderer
295 through a slice cross-attention module. Let \mathbf{F}_b represent
296 the bottleneck feature obtained from the encoder of \mathbf{R} . We
297 utilize the latent codes to query this bottleneck feature, using
298 it as a conditional signal to guide the renderer’s synthesis
299 process. The guiding process can be formulated as:

$$300 \quad \mathbf{F}'_b = \mathbf{F}_b + \mathbf{F}_b \cdot \text{Attn}(q = \theta, k = \mathbf{F}_b, v = \mathbf{F}_b), \quad (12)$$

301 where $\text{Attn}(\cdot)$ denotes the cross-attention operation that fuses
302 the latent codes with the bottleneck feature. Then the refined
303 feature \mathbf{F}'_b is decoded as final images.

304 4.4. Training

305 **GazeGaussian Initialization.** Initialization for the 3D Gaus-
306 sians (3DGS) is crucial for stable optimization. Following
307 Gaussian Head Avatar, we initialize the two-stream Gaus-
308 sians using the neutral mesh extracted from an SDF field.
309 This neutral mesh provides a coarse geometry and texture,
310 which are used to initialize the positions and features of
311 the Gaussians. To decouple the face-only and eye regions,
312 we compute the 3D neutral landmarks and use learnable
313 parameters to define the vertices near the eyes as the initial
314 Gaussians for the eye region, while the rest of the head is
315 used to initialize the face-only Gaussians. Additionally, we
316 transfer the parameters of all deformation and color MLPs
317 while the MLPs for attribute prediction and the expression-
318 guided neural renderer are randomly initialized.

319 **Image Synthesis Loss.** The masked ground truth image \mathcal{I}_{gt}
320 is used to supervise the rendered images $\mathcal{I}_f, \mathcal{I}_e, \mathcal{I}_h$, corre-
321 sponding to the face-only, eyes, and head regions, respec-
322 tively. Additionally, we enforce the first three channels of
323 the feature maps $\mathcal{M}_f, \mathcal{M}_e, \mathcal{M}_h$ to learn the RGB colors.
324 For each rendered image and its corresponding feature map,
325 we apply the same loss functions. Taking the rendered eye
326 image as an example, we mask the ground truth image using
327 an eye mask and then apply L1 loss, SSIM loss, and LPIPS
328 loss on the masked image:

$$329 \quad \mathcal{L}_{\mathcal{I}}^e = \|\mathcal{I}_{gt} - \mathcal{I}_e\|_1 + \lambda_{SSIM}(1 - SSIM(\mathcal{I}_{gt}, \mathcal{I}_e)) \quad (13)$$

$$+ \lambda_{VGG} VGG(\mathcal{I}_{gt}, \mathcal{I}_e),$$

330 where $\lambda_{SSIM} = \lambda_{VGG} = 0.1$ is the weight of loss. The
331 image synthesis loss is the sum of the three rendered images
332 and three feature maps:

$$333 \quad \mathcal{L}_{\mathcal{I}} = \mathcal{L}_{\mathcal{I}}^f + \mathcal{L}_{\mathcal{I}}^e + \mathcal{L}_{\mathcal{I}}^h + \mathcal{L}_{\mathcal{M}}^f + \mathcal{L}_{\mathcal{M}}^e + \mathcal{L}_{\mathcal{M}}^h, \quad (14)$$

where $\mathcal{L}_{\mathcal{I}}^f, \mathcal{L}_{\mathcal{I}}^e, \mathcal{L}_{\mathcal{I}}^h$ represent the losses for the rendered face-
only, eye, and head images, respectively. $\mathcal{L}_{\mathcal{M}}^f, \mathcal{L}_{\mathcal{M}}^e, \mathcal{L}_{\mathcal{M}}^h$
represent the losses for the feature maps corresponding to
the face-only, eye, and head regions, respectively. The image
synthesis loss ensures the full disentanglement of the eye
and the rest of the face.

Gaze Redirection Loss. To improve task-specific perfor-
mance and eliminate task-relevant inconsistencies between
the target image \mathcal{I}_{gt} and the reconstructed head image \mathcal{I}_h ,
we adopt the functional loss used in STED [58] and GazeN-
eRF [34]. The gaze redirection loss can be formulated as:

$$334 \quad \mathcal{L}_G(\mathcal{I}_h, \mathcal{I}_{gt}) = \mathcal{E}_{\text{ang}}(\psi^g(\mathcal{I}_{wf}), \psi^g(\mathcal{I}_{gt})) \quad (15)$$

$$335 \quad \mathcal{E}_{\text{ang}}(\mathbf{v}, \hat{\mathbf{v}}) = \arccos \frac{\mathbf{v} \cdot \hat{\mathbf{v}}}{\|\mathbf{v}\| \|\hat{\mathbf{v}}\|}, \quad 336$$

where $\psi^g(\cdot)$ represents the gaze direction estimated by a
pre-trained gaze estimator network, and $\mathcal{E}_{\text{ang}}(\cdot, \cdot)$ represents
the angular error function. Our final loss function is:

$$337 \quad \mathcal{L} = \lambda_{\mathcal{I}} \mathcal{L}_{\mathcal{I}} + \lambda_G \mathcal{L}_G, \quad (16) \quad 338$$

where $\lambda_{\mathcal{I}} = 1.0$ and $\lambda_G = 0.1$. GazeGaussian is trained
with the final loss until convergence. 339

340 5. Experiments 341

To demonstrate the effectiveness of GazeGaussian, we first
conduct a within-dataset comparison on the ETH-XGaze
dataset [57], testing GazeGaussian alongside state-of-the-art
gaze redirection and head generation methods. Next, we
perform a cross-dataset comparison on ColumbiaGaze [37],
MPIIFaceGaze [54, 55], and GazeCapture [20] to assess gen-
eralization. We also conduct an ablation study to analyze the
contributions of each component in GazeGaussian. Addi-
tionally, we validate the impact of synthesized data on gaze
estimator performance in the supplementary materials. Due
to space limitations, please refer to the supplementary for
more details on the experiment and visualization results. 342

343 5.1. Experimental Settings 344

Dataset Pre-processing. Following GazeNeRF’s prepro-
cessing, we normalize raw images [38, 56] and resize them
into a resolution 512×512 . To enable separate rendering of
the face and eyes regions, we generate masks using face pars-
ing models [62]. We also use the 3D face tracking method
from [50] to produce identity and expression codes and cam-
era poses for the input of our method. For consistency, gaze
labels are converted to pitch-yaw angles in the head coordi-
nate system across all datasets. Details are provided in the
supplementary materials. 345

Baselines. We compare our method with the self-
supervised gaze redirection approach STED [58], along with 346

Table 1. Within-dataset comparison: Quantitative results of the GazeGaussian to other SOTA methods on the ETH-XGaze dataset in terms of gaze and head redirection errors in degree, rendered image quality (SSIM, PSNR, LPIPS, FID), identity similarity and rendering FPS.

Method	Gaze↓	Head Pose↓	SSIM↑	PSNR↑	LPIPS↓	FID↓	Identity Similarity↑	FPS↑
STED	16.217	13.153	0.726	17.530	0.300	115.020	24.347	18
HeadNeRF	12.117	4.275	0.720	15.298	0.294	69.487	46.126	35
GazeNeRF	6.944	3.470	0.733	15.453	0.291	81.816	45.207	46
Gaussian Head Avatar	30.963	13.563	0.638	12.108	0.359	74.560	27.272	91
GazeGaussian (Ours)	6.622	2.128	0.823	18.734	0.216	41.972	67.749	74

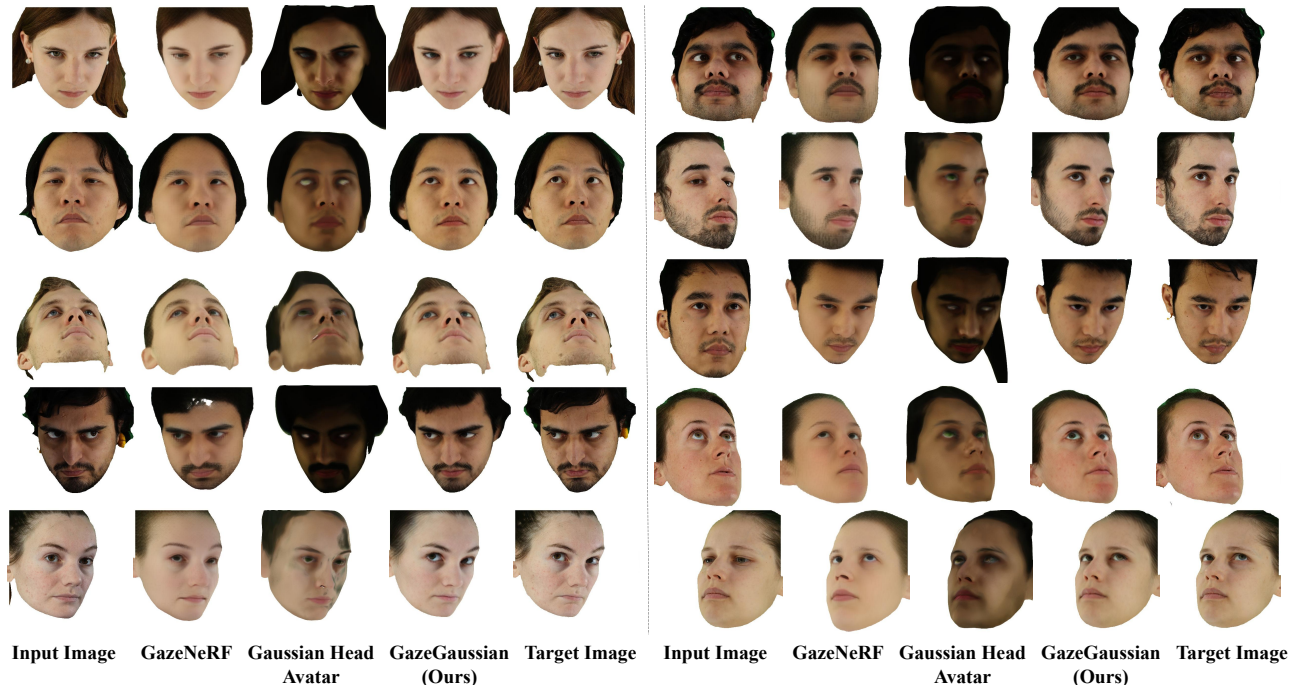


Figure 3. Within-dataset comparison: Visualization of generated images from the ETH-XGaze test set using our GazeGaussian, GazeNeRF, and Gaussian Head Avatar. All faces are masked to remove the background. GazeGaussian generates photo-realistic images with the target gaze direction, preserving identity and facial details. In contrast, GazeNeRF loses identity information and facial details, while Gaussian Head Avatar fails to manipulate the gaze direction effectively.

378 NeRF-based models such as HeadNeRF [15] and the state-of-
 379 the-art method GazeNeRF [34], as well as the latest 3DGS-
 380 based head synthesis method, Gaussian Head Avatar [50].
 381 As the NeRF-based methods, NeRF-Gaze [51] and Wang *et*
 382 *al.* [40] are not yet open-sourced, they are not available for
 383 inclusion in our comparisons.

384 **Metrics.** We evaluate all models using four categories:
 385 redirection accuracy, image quality, identity preservation,
 386 and rendering speed. Redirection accuracy is measured by
 387 gaze and head poses angular errors, using a ResNet50 [13]-
 388 based estimator, as in GazeNeRF [34]. Image quality is
 389 assessed with SSIM, PSNR, LPIPS, and FID. Identity preser-
 390 vation is evaluated with FaceX-Zoo [41], comparing identity
 391 consistency between redirected and ground-truth images.
 392 Rendering speed is reported as average FPS.

5.2. Within-dataset Comparison

393
 394 Following the experimental setup of GazeNeRF, we perform
 395 a within-dataset evaluation to compare the performance of
 396 GazeGaussian with other state-of-the-art methods. All mod-
 397 els are trained using 14.4K images derived from 10 frames
 398 per subject, with 18 camera view images per frame, covering
 399 80 subjects in the ETH-XGaze training set. The evaluation is
 400 conducted on the person-specific test set of the ETH-XGaze
 401 dataset. This test set consists of 15 subjects, each with 200
 402 images annotated with gaze and head pose labels. We follow
 403 the pairing setting in GazeNeRF, which pairs these 200 la-
 404 beled images per subject as input and target samples, and the
 405 same pairings are used across all models to ensure fairness.

406 Tab. 1 presents the quantitative results of GazeGaussian
 407 alongside baseline methods. It can be observed that Gaze-

Table 2. Cross-dataset comparison: Quantitative results of GazeGaussian to other SOTA baselines on ColumbiaGaze, MPIIFaceGaze, and GazeCapture datasets in terms of gaze and head redirection errors in degree, LPIPS, and Identity similarity (ID).

Method	ColumbiaGaze				MPIIFaceGaze				GazeCapture			
	Gaze↓	Head↓	LPIPS↓	ID↑	Gaze↓	Head↓	LPIPS↓	ID↑	Gaze↓	Head↓	LPIPS↓	ID↑
STED	17.887	14.693	0.413	6.384	14.796	11.893	0.288	10.677	15.478	16.533	0.271	6.807
HeadNeRF	15.250	6.255	0.349	23.579	14.320	9.372	0.288	31.877	12.955	10.366	0.232	20.981
GazeNeRF	9.464	3.811	0.352	23.157	14.933	7.118	0.272	30.981	10.463	9.064	0.232	19.025
Gaussian Head Avatar	10.939	3.953	0.347	46.183	12.021	8.530	0.295	30.932	11.571	7.664	0.295	22.236
GazeGaussian (Ours)	7.415	3.332	0.273	59.788	10.943	5.685	0.224	41.505	9.752	7.061	0.209	44.007

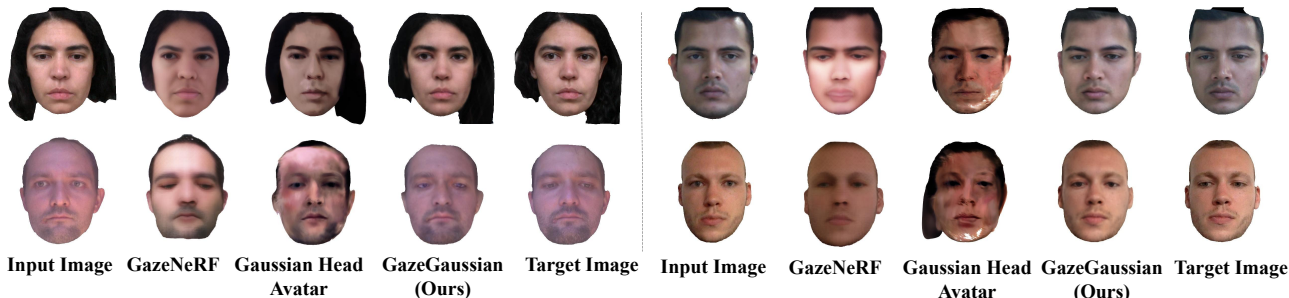


Figure 4. Cross-dataset comparison: Visualization of generated images from the MPIIFaceGaze test set using our GazeGaussian, GazeNeRF, and Gaussian Head Avatar. Please refer to the supplementary for more visualization.

408 Gaussian consistently outperforms prior methods across all
 409 metrics. Specifically, our approach achieves the lowest er-
 410 rors in both gaze and head redirection (6.622° and 2.128° ,
 411 respectively), demonstrating its superior precision in gaze
 412 and head control. Compared to the previous SOTA method
 413 GazeNeRF, which applies rotation to feature map for gaze
 414 redirection, GazeGaussian adopts a Gaussian eye rotation to
 415 explicitly control eye movement. Such a technique not only
 416 improves redirection accuracy but also significantly boosts
 417 rendering quality. Additionally, GazeGaussian achieves a
 418 rendering speed of 74 FPS, nearly doubling the performance
 419 of GazeNeRF, underscoring its efficiency. In contrast, Gaus-
 420 sian Head Avatar (GHA), the latest model built on Gaussian-
 421 based representations, struggles to deliver competitive per-
 422 formance in gaze and head redirection tasks. The lack of
 423 dedicated mechanisms for gaze disentanglement and explicit
 424 eye region modeling in GHA leads to poor performance. By
 425 decoupling the face and eye representation with two-stream
 426 Gaussians, GazeGaussian offers both higher accuracy and
 427 better visual quality, particularly in challenging scenarios
 428 involving extreme head poses or subtle gaze variations.

429 We present a qualitative comparison of different methods
 430 in Fig. 3. GHA struggles to preserve personal identity in the
 431 generated face images, which is quantitatively verified as the
 432 low ‘identity similarity’ in Tab. 1. Moreover, GHA produces
 433 blurred and unrealistic eye regions, significantly degrading
 434 the visual quality of gaze redirection. GazeNeRF, which
 435 implicitly rotates the feature map, fails to effectively control

eye appearance under extreme gaze directions (as shown
 in the last row). Furthermore, it struggles with rendering
 fine-grained facial details and exhibits notable artifacts in
 hair rendering, particularly in the last two rows. Overall, the
 inability to accurately handle eye details in both GHA and
 GazeNeRF limits their effectiveness in gaze redirection. In
 contrast, GazeGaussian consistently produces highly realis-
 tic results, even under challenging conditions, setting a new
 benchmark for gaze redirection tasks.

5.3. Cross-dataset Comparison

To access the generalization capability of GazeGaussian, we
 perform a cross-dataset evaluation on three other datasets:
 ColumbiaGaze, MPIIFaceGaze, and the test set of GazeCap-
 ture. The training setup remains consistent with the within-
 dataset evaluation, using the same model configurations and
 trained parameters.

The results shown in Tab. 2 and Fig. 4 demonstrate that
 GazeGaussian consistently outperforms all other methods
 across the three datasets and all evaluation metrics. By
 introducing a novel expression-guided neural renderer, Gaze-
 Gaussian can retain facial details across various subjects. On
 the other hand, GHA’s performance is limited by its model-
 ing strategy, showing poor adaptability to unseen datasets.
 It produces less clear eye regions and achieves significantly
 lower identity similarity scores compared to GazeGaussian.
 These results further validate the superiority of GazeGaus-
 sian, making it a more robust choice for handling diverse

Table 3. Component-wise ablation study of GazeGaussian on the ETH-XGaze dataset in terms of gaze and head redirection errors in degree, redirection image quality (SSIM, PSNR, LPIPS and FID), and identity similarity.

Two-stream	Gaussian Eye Rep.	Expression-Guided	Gaze↓	Head Pose↓	SSIM↑	PSNR↑	LPIPS↓	FID↓	Identity Similarity↑
✓			13.651	2.981	0.753	16.376	0.272	55.481	38.941
✓		✓	13.489	3.149	0.751	16.365	0.274	54.327	38.521
✓	✓		8.883	2.635	0.766	16.692	0.254	48.891	45.013
	✓	✓	7.494	3.098	0.769	16.873	0.250	49.658	46.155
✓	✓	✓	6.622	2.128	0.823	18.734	0.216	41.972	67.749

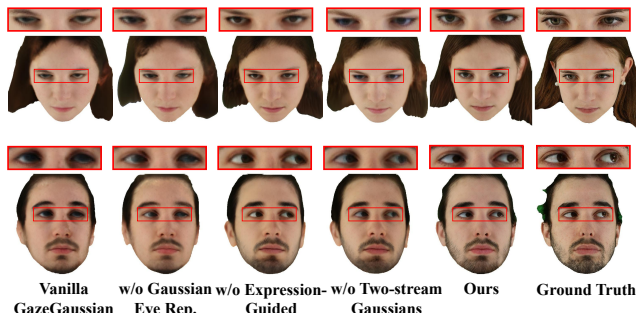


Figure 5. Qualitative ablation study on the ETH-XGaze dataset.

463 datasets and complex gaze redirection tasks. Please refer to
 464 supplementary material for more visualization results on the
 465 cross-dataset evaluation.

466 5.4. Ablation Study

467 To validate the effectiveness of each component, we conduct
 468 a component-wise ablation study on the ETH-XGaze dataset.
 469 The results are shown in Tab. 3 and Fig. 5.

470 **Vanilla-GazeGaussian.** In this version, we omit the pro-
 471 posed Gaussian eye rotation representation and expression-
 472 guided neural renderer. The corresponding experimental
 473 results are shown in the first row of the table and the first
 474 column of the visualizations. The eye deformation is treated the
 475 same as the face, and the neural renderer remains unchanged
 476 from GazeNeRF. The results show that, due to the lack of
 477 control over eye rotation, gaze redirection errors are large,
 478 and the image synthesis quality is relatively low.

479 **w/o Gaussian eye rotation representation.** To verify the
 480 contribution of the proposed Gaussian eye rotation repre-
 481 sentation, we omit it in the GazeGaussian. The results are
 482 shown in the second row of the table and the second column
 483 of the figure. Compared to the full version of GazeGaussian,
 484 the introduction of a specialized representation for eye de-
 485 formation significantly improves gaze redirection accuracy
 486 and enhances the detail in the eye region.

487 **w/o Expression-Guided.** We remove the proposed
 488 expression-guided neural renderer and rely solely on the
 489 neural renderer in GazeNeRF for image synthesis. The
 490 results, shown in the third row of the table and the third

column of the figure, indicate a noticeable decline in image
 quality. Without expression guidance, the model struggles to
 effectively preserve dynamic facial expressions, leading to
 less accurate gaze redirection. The synthesized images also
 exhibit lower fidelity in capturing facial details and subtle
 expression changes.

w/o Two-stream. Replacing the two-stream structure with a
 single-stream Gaussian model for both face and eye regions
 leads to performance degradation and loss of synthesis de-
 tails, as shown in the fourth row of the table and the fourth
 column of the figure. Combining face and eye regions in
 a single stream fails to capture the eye region’s complex
 dynamics, resulting in less accurate gaze redirection and
 lower image fidelity. The two-stream architecture, which
 decouples the face and eye regions, enables more precise
 modeling of each region’s unique characteristics, improv-
 ing gaze accuracy and image quality. Furthermore, when
 comparing this version to the vanilla GazeGaussian (where
 no proposed components are used), we observe a substan-
 tial performance improvement, validating the effectiveness
 of the proposed techniques and their contribution to gaze
 redirection and head avatar synthesis.

Among all the ablation experiments, the full GazeGaus-
 sian achieves the best performance. This improvement re-
 sults from the combination of the two-stream Gaussian struc-
 ture, which decouples the face and eye regions for more
 precise modeling, and the proposed Gaussian eye rotation
 representation, which enables accurate control of eye rota-
 tion. Additionally, the expression-guided neural renderer
 enhances the model’s ability to generalize across subjects
 while preserving facial details.

522 6. Conclusion

523 We present GazeGaussian, a high-fidelity gaze redirection
 524 pipeline that uses a two-stream model to represent face and
 525 eye regions separately. We present a new Gaussian-based
 526 representation of the eye to accurately depict eye rotations,
 527 along with an expression-conditional neural renderer that
 528 enhances the fidelity of gaze redirection. Numerous experi-
 529 ments have shown that GazeGaussian achieves state-of-the-
 530 art performance on the task of gaze direction, paving the way
 531 for more robust gaze estimation in real-world applications.

532

533

References

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

ings of the IEEE conference on computer vision and pattern recognition, pages 770–778, 2016. 6

[14] Zhe He, Adrian Spurr, Xucong Zhang, and Otmar Hilliges. Photo-realistic monocular gaze redirection using generative adversarial networks. In *IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2019. 2

[15] Yang Hong, Bo Peng, Haiyao Xiao, Ligang Liu, and Juyong Zhang. Headnerf: A real-time nerf-based parametric head model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20374–20384, 2022. 6

[16] Yang Hong, Bo Peng, Haiyao Xiao, Ligang Liu, and Juyong Zhang. Headnerf: A real-time nerf-based parametric head model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20374–20384, 2022. 1, 2

[17] Nan Huang, Xiaobao Wei, Wenzhao Zheng, Pengju An, Ming Lu, Wei Zhan, Masayoshi Tomizuka, Kurt Keutzer, and Shanghang Zhang. S3gaussian: Self-supervised street gaussians for autonomous driving. *arXiv preprint arXiv:2405.20323*, 2024. 1

[18] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 1, 2, 3

[19] Daniil Kononenko and Victor Lempitsky. Learning to look up: Realtime monocular gaze correction using machine learning. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4667–4675, 2015. 2

[20] Kyle Krafka, Aditya Khosla, Petr Kellnhofer, Harini Kannan, Suchendra Bhandarkar, Wojciech Matusik, and Antonio Torralba. Eye tracking for everyone. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2176–2184, 2016. 5

[21] Claudia Kuster, Tiberiu Popa, Jean-Charles Bazin, Craig Gotsman, and Markus Gross. Gaze correction for home video conferencing. *ACM Trans. Graph.*, 31(6), 2012. 2

[22] Gengyan Li, Abhimitra Meka, Franziska Mueller, Marcel C Buehler, Otmar Hilliges, and Thabo Beeler. Eyenerf: a hybrid representation for photorealistic synthesis, animation and relighting of human eyes. *ACM Transactions on Graphics (TOG)*, 41(4):1–16, 2022. 2

[23] Tianye Li, Timo Bolkart, Michael J Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4d scans. *ACM Trans. Graph.*, 36(6):194–1, 2017. 2

[24] Tao Lu, Mulin Yu, Linning Xu, Yuanbo Xiangli, Limin Wang, Dahua Lin, and Bo Dai. Scaffold-gs: Structured 3d gaussians for view-adaptive rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20654–20664, 2024. 1

[25] Annu George Mavelly, JE Judith, PA Sahal, and Steffy Ann Kuruvilla. Eye gaze tracking based driver monitoring system. In *2017 IEEE international conference on circuits and systems (ICCS)*, pages 364–367. IEEE, 2017. 1

[26] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis, 2020. 1, 2

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

- 646 [27] Nitish Padmanaban, Robert Konrad, Emily A Cooper, and
647 Gordon Wetzstein. Optimizing vr for all users through adap-
648 tive focus displays. In *ACM SIGGRAPH 2017 Talks*, pages
649 1–2, 2017. 1
- 650 [28] Seonwook Park, Shalini De Mello, Pavlo Molchanov, Umar
651 Iqbal, Otmar Hilliges, and Jan Kautz. Few-shot adaptive gaze
652 estimation. In *Proceedings of the IEEE/CVF international
653 conference on computer vision*, pages 9368–9377, 2019. 2
- 654 [29] Ziqiao Peng, Yihao Luo, Yue Shi, Hao Xu, Xiangyu Zhu,
655 Hongyan Liu, Jun He, and Zhaoxin Fan. Selftalk: A self-
656 supervised commutative training diagram to comprehend 3d
657 talking faces. In *Proceedings of the 31st ACM International
658 Conference on Multimedia*, pages 5292–5301, 2023. 2
- 659 [30] Ziqiao Peng, Haoyu Wu, Zhenbo Song, Hao Xu, Xiangyu
660 Zhu, Jun He, Hongyan Liu, and Zhaoxin Fan. Emotalk:
661 Speech-driven emotional disentanglement for 3d face anima-
662 tion. In *Proceedings of the IEEE/CVF International Confer-
663 ence on Computer Vision*, pages 20687–20697, 2023. 2
- 664 [31] Shenhan Qian, Tobias Kirschstein, Liam Schoneveld, Davide
665 Davoli, Simon Giebenhain, and Matthias Nießner. Gaussian-
666 avatars: Photorealistic head avatars with rigged 3d gaussians.
667 In *Proceedings of the IEEE/CVF Conference on Computer
668 Vision and Pattern Recognition*, pages 20299–20309, 2024.
669 1, 2
- 670 [32] Yalun Qin, Kuo-Chin Lien, Matthew Turk, and Tobias
671 Höllerer. Eye gaze correction with a single webcam based on
672 eye-replacement. In *Advances in Visual Computing*, pages
673 599–609, Cham, 2015. Springer International Publishing. 2
- 674 [33] Anurag Ranjan, Timo Bolkart, Soubhik Sanyal, and Michael J
675 Black. Generating 3d faces using convolutional mesh au-
676 toencoders. In *Proceedings of the European conference on
677 computer vision (ECCV)*, pages 704–720, 2018. 2
- 678 [34] Alessandro Ruzzi, Xiangwei Shi, Xi Wang, Gengyan Li,
679 Shalini De Mello, Hyung Jin Chang, Xucong Zhang, and
680 Otmar Hilliges. Gazenerf: 3d-aware gaze redirection with
681 neural radiance fields. In *2023 IEEE/CVF Conference on
682 Computer Vision and Pattern Recognition (CVPR)*, pages
683 9676–9685, 2023. 1, 2, 3, 5, 6
- 684 [35] Tianchang Shen, Jun Gao, Kangxue Yin, Ming-Yu Liu, and
685 Sanja Fidler. Deep marching tetrahedra: a hybrid represen-
686 tation for high-resolution 3d shape synthesis. *Advances in
687 Neural Information Processing Systems*, 34:6087–6101, 2021.
688 3
- 689 [36] Zhixin Shu, Eli Shechtman, Dimitris Samaras, and Sunil
690 Hadap. Eyeopener: Editing eyes in the wild. *ACM Trans.
691 Graph.*, 36(1), 2016. 2
- 692 [37] B.A. Smith, Q. Yin, S.K. Feiner, and S.K. Nayar. Gaze
693 Locking: Passive Eye Contact Detection for Human-Object
694 Interaction. In *ACM Symposium on User Interface Software
695 and Technology (UIST)*, pages 271–280, 2013. 5
- 696 [38] Yusuke Sugano, Yasuyuki Matsushita, and Yoichi Sato.
697 Learning-by-synthesis for appearance-based 3d gaze estima-
698 tion. In *Proceedings of the IEEE conference on computer
699 vision and pattern recognition*, pages 1821–1828, 2014. 5
- 700 [39] Luan Tran and Xiaoming Liu. Nonlinear 3d face morphable
701 model. In *Proceedings of the IEEE conference on computer
702 vision and pattern recognition*, pages 7346–7355, 2018. 1
- [40] Hengfei Wang, Zhongqun Zhang, Yihua Cheng, and 703
Hyung Jin Chang. High-fidelity eye animatable neural radi- 704
ance fields for human face. *arXiv preprint arXiv:2308.00773*, 705
2023. 1, 2, 6 706
- [41] Jun Wang, Yinglu Liu, Yibo Hu, Hailin Shi, and Tao Mei. 707
Facex-zoo: A pytorch toolbox for face recognition. In *Pro- 708
ceedings of the 29th ACM International Conference on Multi- 709
media*, pages 3779–3782, 2021. 6 710
- [42] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku 711
Komura, and Wenping Wang. Neus: Learning neural implicit 712
surfaces by volume rendering for multi-view reconstruction. 713
arXiv preprint arXiv:2106.10689, 2021. 1 714
- [43] Yu Wang, Xiaobao Wei, Ming Lu, and Guoliang Kang. Plgs: 715
Robust panoptic lifting with 3d gaussian splatting. *arXiv 716
preprint arXiv:2410.17505*, 2024. 1 717
- [44] Xiaobao Wei, Renrui Zhang, Jiarui Wu, Jiaming Liu, Ming 718
Lu, Yandong Guo, and Shanghang Zhang. Nto3d: Neural 719
target object 3d reconstruction with segment anything. In *Pro- 720
ceedings of the IEEE/CVF Conference on Computer Vision 721
and Pattern Recognition*, pages 20352–20362, 2024. 1 722
- [45] Erroll Wood, Tadas Baltrusaitis, Louis-Philippe Morency, 723
Peter Robinson, and Andreas Bulling. Gazedirector: Fully 724
articulated eye gaze redirection in video, 2017. 2 725
- [46] Weihao Xia, Yujtu Yang, Jing-Hao Xue, and Wensen Feng. 726
Controllable continuous gaze redirection. In *Proceedings of 727
the 28th ACM International Conference on Multimedia*, pages 728
1782–1790, 2020. 2 729
- [47] Jun Xiang, Xuan Gao, Yudong Guo, and Juyong Zhang. 730
Flashavatar: High-fidelity head avatar with efficient gaussian 731
embedding. In *The IEEE Conference on Computer Vision 732
and Pattern Recognition (CVPR)*, 2024. 1, 2 733
- [48] Mingjie Xu, Haofei Wang, and Feng Lu. Learning a general- 734
ized gaze estimator from gaze-consistent feature. In *Proce- 735
edings of the AAAI conference on artificial intelligence*, pages 736
3027–3035, 2023. 1 737
- [49] Yuelang Xu, Benwang Chen, Zhe Li, Hongwen Zhang, 738
Lizhen Wang, Zerong Zheng, and Yebin Liu. Gaussian head 739
avatar: Ultra high-fidelity head avatar via dynamic gaussians. 740
arXiv preprint arXiv:2312.03029, 2023. 2 741
- [50] Yuelang Xu, Benwang Chen, Zhe Li, Hongwen Zhang, 742
Lizhen Wang, Zerong Zheng, and Yebin Liu. Gaussian head 743
avatar: Ultra high-fidelity head avatar via dynamic gaussians. 744
In *Proceedings of the IEEE/CVF Conference on Computer 745
Vision and Pattern Recognition (CVPR)*, 2024. 1, 2, 3, 5, 6 746
- [51] Pengwei Yin, Jingjing Wang, Jiawu Dai, and Xiaojun Wu. 747
Nerf-gaze: A head-eye redirection parametric model for gaze 748
estimation. In *ICASSP 2024-2024 IEEE International Confer- 749
ence on Acoustics, Speech and Signal Processing (ICASSP)*, 750
pages 2760–2764. IEEE, 2024. 1, 2, 6 751
- [52] Yu Yu and Jean-Marc Odobez. Unsupervised representation 752
learning for gaze estimation. In *Proceedings of the IEEE/CVF 753
Conference on Computer Vision and Pattern Recognition*, 754
pages 7314–7324, 2020. 1 755
- [53] Mingfang Zhang, Yunfei Liu, and Feng Lu. Gazeonce: 756
Real-time multi-person gaze estimation. In *Proceedings of 757
the IEEE/CVF Conference on Computer Vision and Pattern 758
Recognition*, pages 4197–4206, 2022. 1 759

- 760 [54] Xucong Zhang, Yusuke Sugano, Mario Fritz, and Andreas
761 Bulling. Appearance-based gaze estimation in the wild. In
762 *Proc. of the IEEE Conference on Computer Vision and Pattern
763 Recognition (CVPR)*, pages 4511–4520, 2015. 5
- 764 [55] Xucong Zhang, Yusuke Sugano, Mario Fritz, and Andreas
765 Bulling. It’s written all over your face: Full-face appearance-
766 based gaze estimation. In *Computer Vision and Pattern Recog-
767 nition Workshops (CVPRW), 2017 IEEE Conference on*, pages
768 2299–2308. IEEE, 2017. 5
- 769 [56] Xucong Zhang, Yusuke Sugano, and Andreas Bulling. Revis-
770 iting data normalization for appearance-based gaze estimation.
771 In *Proc. International Symposium on Eye Tracking Research
772 and Applications (ETRA)*, pages 12:1–12:9, 2018. 1, 5
- 773 [57] Xucong Zhang, Seonwook Park, Thabo Beeler, Derek
774 Bradley, Siyu Tang, and Otmar Hilliges. Eth-xgaze: A large
775 scale dataset for gaze estimation under extreme head pose and
776 gaze variation. In *European Conference on Computer Vision
777 (ECCV)*, 2020. 5
- 778 [58] Yufeng Zheng, Seonwook Park, Xucong Zhang, Shalini De
779 Mello, and Otmar Hilliges. Self-learning transformations for
780 improving gaze and head redirection. In *Neural Information
781 Processing Systems (NeurIPS)*, 2020. 2, 5
- 782 [59] Yufeng Zheng, Wang Yifan, Gordon Wetzstein, Michael J
783 Black, and Otmar Hilliges. Pointavatar: Deformable point-
784 based head avatars from videos. In *Proceedings of the
785 IEEE/CVF conference on computer vision and pattern recog-
786 nition*, pages 21057–21067, 2023. 1, 2
- 787 [60] Wojciech Zielonka, Timo Bolkart, and Justus Thies. Towards
788 metrical reconstruction of human faces. In *European con-
789 ference on computer vision*, pages 250–269. Springer, 2022.
790 1
- 791 [61] Wojciech Zielonka, Timo Bolkart, and Justus Thies. Instant
792 volumetric head avatars. In *Proceedings of the IEEE/CVF
793 Conference on Computer Vision and Pattern Recognition*,
794 pages 4574–4584, 2023. 1, 2
- 795 [62] zllrunning. Using modified bisenet for face parsing in pytorch,
796 2019. [https://github.com/zllrunning/face-
797 parsing.PyTorch](https://github.com/zllrunning/face-parsing.PyTorch). 5