

3D Gaussian and Diffusion-Based Gaze Redirection

Abiram Panchalingam

ap18g21@soton.ac.uk

Indu P. Bodala

I.P.Bodala@soton.ac.uk

Stuart E. Middleton

sem03@soton.ac.uk

School of Electronics and Computer Science, University of Southampton

Abstract

High-fidelity gaze redirection is critical for generating augmented data to improve the generalization of gaze estimators. 3D Gaussian Splatting (3DGS) models like Gaze-Gaussian represent the state-of-the-art but can struggle with rendering subtle, continuous gaze shifts. In this paper, we propose DiT-Gaze, a framework that enhances 3D gaze redirection models using a novel combination of Diffusion Transformer (DiT), weak supervision across gaze angles, and an orthogonality constraint loss. DiT allows higher-fidelity image synthesis, while our weak supervision strategy using synthetically generated intermediate gaze angles provides a smooth manifold of gaze directions during training. The orthogonality constraint loss mathematically enforces the disentanglement of internal representations for gaze, head pose, and expression. Comprehensive experiments show that DiT-Gaze sets a new state-of-the-art in both perceptual quality and redirection accuracy, reducing the state-of-the-art gaze error by 4.1% to 6.353 degrees, providing a superior method for creating synthetic training data. Our code and models will be made available for the research community to benchmark against.

1. Introduction

Gaze redirection is a critical task for augmenting datasets to improve the generalization of gaze estimators, which often fail in cross-domain scenarios when encountering out-of-distribution data. This process, which involves manipulating a person's gaze in an image while preserving their identity, generates synthetic data that can enhance the robustness of these estimators. Early works in this domain formulated gaze redirection as a 2D image-to-image translation problem, relying on techniques like image warping [3] or generative models such as Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs) [4, 30]. However, these 2D methods fundamentally overlooked the 3D nature of head and eye movements, often resulting in poor spatial consistency, visual artifacts, and a limited range of redirection.

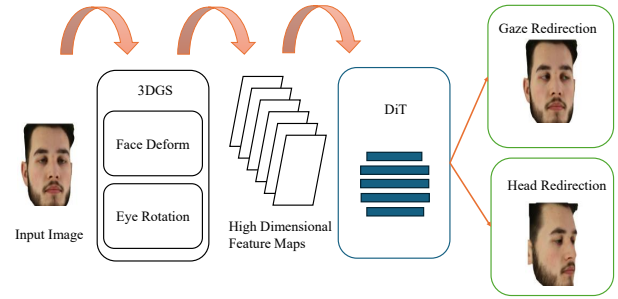


Figure 1. Gaze redirection: Given an input image and target gaze, DiTGaze utilizes a 3DGS model and a DiT renderer to generate high-fidelity head images with accurate gaze redirection.

To address these limitations, the task was reformulated as a 3D-aware problem, leveraging advancements in neural rendering. **GazeNeRF** [19] was a pioneering method that used **Neural Radiance Fields (NeRF)** [11] to model the face. It introduced a **two-stream MLP architecture** to separately represent the face and eye regions, enabling a 3D rotation to be applied to the volumetric eye features before rendering. While this approach significantly improved 3D consistency, GazeNeRF and other NeRF-based methods are hindered by the high computational demands and slow rendering speeds inherent to volumetric rendering.

The current state-of-the-art, **GazeGaussian** [24], advanced the field by being the first to leverage **3D Gaussian Splatting (3DGS)** [8] for this task, overcoming the speed limitations of NeRF. By adopting a similar two-stream model and introducing a novel eye rotation field for explicit control, it achieves higher redirection accuracy and significantly faster rendering speeds than its NeRF-based predecessors. However, despite its success, its U-Net-based renderer presents an opportunity for improvement with a more powerful generative architecture, and like many models, it can struggle to synthesize the subtle, intermediate gaze angles not explicitly seen in the discrete training pairs.

In this paper, we propose DiT-Gaze, a 3D Gaussian Splatting-based framework with two key contributions. First, we use the Diffusion Transformer (DiT) architecture

for the first time in the field of gaze redirection, leveraging the superior scalability and self-attention mechanism of transformers for higher-fidelity image synthesis [16]. In addition, to render a smooth manifold of gaze redirection, we introduce a weak supervision strategy directly in the 3D Gaussian space. Second, we implement a novel Orthogonality Constraint Loss to mathematically enforce the disentanglement of the internal representations for gaze, head pose, and expression, a technique not previously applied to 3DGS-based avatars.

2. Related Work

2.1. Gaze Redirection

Early approaches formulated gaze redirection as a 2D image manipulation task [3, 22]. These methods included image warping techniques such as DeepWarp [3], which were limited in their ability to handle large gaze shifts, and generative models like conditional GANs [2, 4], which improved image quality but fundamentally lacked 3D awareness. Other methods such as STED [30] and ReDirTrans [7] also operated in 2D by applying rotations in a learned latent space, often resulting in poor spatial consistency and visual artifacts.

The limitations of 2D manipulation prompted a shift toward 3D-aware solutions. GazeNeRF [19] led this transition by employing Neural Radiance Fields (NeRF). It introduced a critical architectural principle: a two-stream model that **decouples** the representation of the static face from the dynamic eye regions. This separation was designed to achieve disentangled control, allowing the gaze to be manipulated independently of the head pose. In this framework, gaze redirection is performed by applying a rigid 3D rotation to the eye region’s neural field before compositing it with the face field. Despite its improved consistency, the method’s reliance on NeRF’s volumetric rendering resulted in significant computational overhead and slow inference. Moreover, this feature-level manipulation can still struggle to preserve identity and fine-grained textures, often producing soft or blurred results.

GazeGaussian [24] represents the state-of-the-art by being the first to leverage 3D Gaussian Splatting (3DGS) [8] for this task, overcoming the speed limitations of NeRF. Adopting a similar two-stream model, GazeGaussian introduces a novel **Eye Rotation field** that explicitly adjusts the positions of the 3D eye Gaussians to simulate a rigid rotation. This explicit control mechanism, combined with the efficiency of 3DGS, allows it to achieve higher redirection accuracy and significantly faster rendering speeds than its NeRF-based predecessors.

2.2. 3D Head Avatar Synthesis

The synthesis of dynamic 3D head avatars is a foundational area of research that provides the underlying technology for 3D-aware gaze redirection. Early work in this domain utilized parametric 3D head models like FLAME [1] to map expression and pose parameters directly to 3D facial geometry. More recent work has shifted to neural rendering techniques, which can be broadly categorized into NeRF-based and 3DGS-based methods.

One major approach uses **Neural Radiance Fields (NeRF)** to create high-fidelity, controllable avatars. Models such as **HeadNeRF** [6] leverage neural radiance fields to deform facial movements from a canonical space, conditioned on parameters for shape, expression, and lighting. However, these methods are often limited by the high computational demands and slow rendering speeds inherent to volumetric rendering.

More recently, **3D Gaussian Splatting (3DGS)** has emerged as a superior alternative, offering impressive rendering quality at significantly faster speeds. State-of-the-art models like **Gaussian Head Avatar** [26] initialize 3D Gaussians from a neutral mesh and use MLPs to deform them, creating ultra high-fidelity dynamic avatars. While these general head avatar models are powerful for facial animation, they often neglect the mechanisms for precise gaze control, a key limitation that specialized models like GazeGaussian are designed to address.

2.3. Diffusion Models for Image Synthesis

Diffusion models have become the state-of-the-art for high-fidelity image synthesis, operating through a forward process that gradually adds noise to data and a learned reverse process that iteratively denoises it to generate a clean sample [5].

Early diffusion models often operated directly in the pixel space, using convolutional U-Net backbones [18] to predict the noise component at each step of the reverse process. To improve computational efficiency, Latent Diffusion Models (LDMs) [17] were introduced, which operate in a compressed latent space created by a pre-trained autoencoder, significantly reducing complexity without sacrificing quality.

The most recent advancement in this domain is the Diffusion Transformer (DiT) [16], introduced by Peebles and Xie, which replaces the U-Net backbone with a more scalable and powerful transformer architecture. This trend of replacing convolutional U-Nets with transformers has been validated in various specialized domains, including 3D shape generation [12], medical image synthesis [14], and scientific climate simulation [13]. DiTs operate on a sequence of latent patches and have demonstrated superior performance and scaling properties in image generation benchmarks. The key advantage of a DiT is its ability

to model long-range dependencies through self-attention, leading to better global coherence and higher-quality image synthesis compared to the more locally focused convolutions of a U-Net. This proven superiority in generative capability and scalability motivates our work to replace GazeGaussian’s U-Net renderer with a DiT for improved fidelity.

3. Method

3.1. Framework Overview

Our framework enhances the GazeGaussian architecture [24], a two-stream 3DGS model that initializes face and eye Gaussians from a neutral mesh [26]. This baseline deforms the canonical Gaussians using pose and expression, rotates them via gaze direction, and rasters them into a feature map.

We introduce three novel contributions. First, we replace the baseline’s renderer with a **Diffusion Transformer (DiT) Neural Renderer**, which uses **AdaLN (Adaptive Layer Norm)** [16] for conditioning. Second, we implement a **weak supervision strategy**, training on intermediate gaze angles to create a smooth gaze manifold. Third, we add a novel **Orthogonality Constraint Loss** to enforce feature disentanglement by reusing existing MLP features.

3.2. Two-Stream 3D Gaussian Representation

GazeGaussian’s architecture is based on a two-stream 3DGS model that decouples the face and eye regions into two separate, independently controlled sets of 3D Gaussians. The pipeline begins by learning an implicit Signed Distance Function (SDF) from the training data, from which a neutral 3D mesh is extracted using DMTet [21], a process adapted from Gaussian Head Avatar [25]. This mesh is then partitioned using 3D landmarks to initialize the canonical Gaussians for the face-only stream and the eye stream, providing a robust starting point for deformation and rotation.

The Face Deformation Field is responsible for manipulating the face-only stream. It begins with a set of canonical neutral face Gaussians, each with attributes for position, features, rotation, scale, and opacity:

$$\{\mu_0^f, z_0^f, R_0^f, S_0^f, \alpha_0^f\} \quad (1)$$

A set of MLPs, conditioned on head pose (γ) and expression (τ) codes, predicts the final transformed state of these Gaussians. The influence of pose and expression is blended using learned weights (λ_τ and λ_γ), which are determined by each Gaussian’s proximity to 3D facial landmarks. The final transformed attributes are calculated as follows:

$$\mu^f = \mu_0^f + \lambda_\tau E_\mu^f(\mu_0^f, \tau) + \lambda_\gamma P_\mu^f(\mu_0^f, \gamma) \quad (2)$$

$$c^f = \lambda_\tau E_c^f(z_0^f, \tau) + \lambda_\gamma P_c^f(z_0^f, \gamma) \quad (3)$$

$$\kappa^f = \kappa_0^f + \lambda_\tau E_\kappa^f(z_0^f, \tau) + \lambda_\gamma P_\kappa^f(z_0^f, \gamma) \quad (4)$$

where κ^f represents the rotation, scale, and opacity attributes. These equations calculate the final transformed attributes for each face Gaussian, blending the outputs of expression-conditioned MLPs (E^f) and pose-conditioned MLPs (P^f), with influence controlled by λ_τ and λ_γ .

Eye Rotation Field is dedicated to the eye stream and utilizes the **Gaussian Eye Rotation Representation** to achieve precise gaze control. This component begins with a set of canonical neutral eye Gaussians:

$$\{\mu_0^e, z_0^e, R_0^e, S_0^e, \alpha_0^e\} \quad (5)$$

where each attribute left to right represents position, features, rotation, scale, and opacity. However, in the Face Deformation Field, the scale parameter is a 3D vector allowing non-uniform scaling, whereas in the Eye Rotation Field, it is constrained to be a single scalar value to maintain a uniform spherical shape that better aligns with the rotational properties of an eyeball.

To simulate rigid eyeball rotation, a separate set of MLPs explicitly computes transformations based on the target gaze vector (ϕ) and expression codes (τ). This allows the model to directly adjust the positions and attributes of the eye Gaussians to align with the desired gaze direction, as shown in the following equations:

$$\mu^e = E_\mu^e(\mu_0^e, \tau) + G_\mu^e(\mu_0^e, \phi) \quad (6)$$

$$c^e = E_c^e(z_0^e, \tau) + G_c^e(z_0^e, \phi) \quad (7)$$

$$\kappa^e = \kappa_0^e + E_\kappa^e(z_0^e, \tau) + G_\kappa^e(z_0^e, \phi) \quad (8)$$

These equations detail how the final attributes of the eye Gaussians are calculated. The final position (μ^e) is determined by combining a deformation offset from an expression-conditioned MLP (E_μ^e) with a rotational transformation from a gaze-conditioned MLP (G_μ^e). Similarly, the final color (c^e) and other attributes (κ^e) are calculated by summing the outputs of separate MLPs conditioned on both the expression code (τ) and the target gaze vector (ϕ).

This explicit, disentangled control over the face and eye streams is the key to GazeGaussian’s high-fidelity redirection capabilities.

3.3. Diffusion Transformer for High-Fidelity Rendering

Our primary architectural contribution is the replacement of GazeGaussian’s U-Net-based Expression-Guided Neural Renderer (EGNR) with a more powerful **Diffusion Transformer (DiT)** architecture, as introduced by Peebles et al. [21]. The DiT can be designed to operate on the rasterized feature map from the 3DGS model, but at a higher cost of performance.

To improve computational efficiency, our DiT operates in a compressed latent space, following the Latent Diffusion

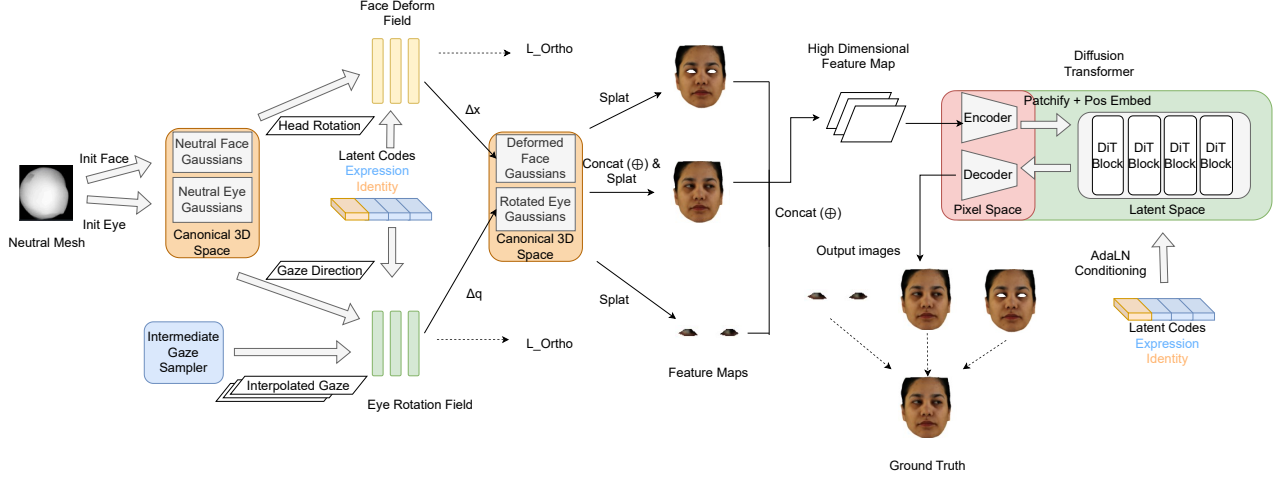


Figure 2. Pipeline of DiTGaze. We initialize face and eye Gaussians from a pre-trained neutral mesh. The Intermediate Gaze Sampler (Ours) generates inputs for the Eye Rotation Field (Δq), while pose and expression codes drive the Face Deform Field (Δx). To ensure disentanglement, our Orthogonality Loss (L_{Ortho}) is applied to both fields. The resulting Gaussians are splatted into feature maps, concatenated, and fed into our Latent DiT Renderer (Ours), which uses AdaLN conditioning to synthesize the final image.

Model framework [17]. The **input processing** begins with the high-dimensional feature map from the 3DGS rasterizer, which is first passed through the pre-trained and frozen VAE encoder [17]. The **PatchEmbed** module uses a single Conv2D layer to efficiently divide the latent map into a grid of patches and linearly project them into a sequence of tokens. To retain spatial information, we add 2D sinusoidal positional embeddings (**pos_embed**) to these tokens before they are fed into the transformer blocks.

The **iterative denoising process** is guided by a combined conditioning vector. The scalar diffusion timestep (τ) is first projected into a high-dimensional vector using a **TimestepEmbedding** module. This embedding is then concatenated with the **shape_code**, which contains the expression and identity information. This combined conditioning vector is then used to modulate the behavior of each **DiTBlock** via an Adaptive Layer Normalization (AdaLN) mechanism, allowing the model to adapt its processing based on both the noise level and the desired facial appearance.

AdaLN was found to be the most effective conditioning strategy in the original DiT paper. The operation can be described as:

$$\text{adaLN}(x, c) = (1 + \text{scale}) \cdot \text{LayerNorm}(x) + \text{shift} \quad (9)$$

where the **scale** and **shift** parameters are regressed from the combined conditioning vector c . This allows the model to adapt its processing based on both the noise level and the desired facial appearance.

The **output generation** is performed by a series of **DiT-Block** modules, which process the sequence of image tokens conditioned on the timestep and expression codes.

After processing, a **FinalLayer** projects the output tokens back into the patch dimension. Finally, an **unpatchify** operation reconstructs the processed latent map, which is then passed through the pre-trained VAE decoder [17] to synthesize the final, high-fidelity, gaze-redirected image. The entire renderer is trained end-to-end as part of the Gaze-Gaussian pipeline to transform the input feature map from the 3DGS rasterizer into the final rendered image.

3.4. Weak Supervision with Gaze Interpolation

To address the challenge of rendering subtle gaze shifts not present in the discrete training pairs, we implement a weak supervision strategy inspired by prior work on fine-grained gaze learning [15]. Rather than relying solely on the sparse target angles from the dataset, we generate a rich distribution of synthetic intermediate gaze vectors on-the-fly during each training step. Our implementation uses a configurable sampling method with three distinct modes: a **uniform** distribution across the gaze range, a **biased_center** distribution that focuses on more stable, central gazes, and a **mixed** mode that combines systematic grid sampling with random sampling to ensure comprehensive coverage. The grid sampling component deterministically generates intermediate angles ($G_{i,j}$) as follows:

$$g_{yaw} = -R_{gaze} + \frac{2 \cdot R_{gaze} \cdot i}{N_{grid} - 1} \quad (10)$$

$$g_{pitch} = -R_{gaze} + \frac{2 \cdot R_{gaze} \cdot j}{N_{grid} - 1} \quad (11)$$

where R_{gaze} is the maximum gaze range, N_{grid} is the grid size, which is the number of discrete points sampled

along each axis (horizontal/yaw and vertical/pitch) of the 2D gaze space, and i, j are indices from 0 to $N_{grid} - 1$.

The purpose of this grid is to systematically ensure that the synthetic gaze angles are evenly distributed across the entire gaze range, from $-R_{gaze}$ to $+R_{gaze}$. This guarantees that the model is trained on a variety of intermediate, in-between angles, preventing gaps in its learned understanding of the continuous gaze space.

Furthermore, to enhance training stability, our strategy is progressive. In early epochs, the model is trained on a smaller, center-biased range of angles. As training progresses, the range is gradually expanded, and the sampling strategy shifts to more uniform distributions to cover more extreme gaze directions. These synthetic gaze vectors are fed as the target input to the **Eye Rotation Field**, which transforms the 3D eye Gaussians accordingly.

The complete model, including the DiT renderer, is then trained to generate the corresponding image for this intermediate state, which is supervised using the standard **Gaze Redirection Loss**. This process explicitly teaches the model to render a smooth and continuous manifold of gaze directions, making the manipulation more interpretable and geometrically grounded in the explicit 3D Gaussian space.

3.5. Orthogonality Constraint Loss

Our third contribution is a novel loss function that mathematically enforces feature disentanglement as a refinement to the model’s structural two-stream design.

We introduce an **Orthogonality Constraint Loss** that penalizes the correlation between the internal representations of gaze, head pose, and expression. To implement this in a computationally efficient manner, we reuse the existing architecture instead of adding new encoder networks. We capture the first internal representation of each control vector by taking the output of the first linear layer of its corresponding MLP within the Gaussian Model:

Vector	MLP
Gaze (φ)	eye_deform_mlp
Pose (γ)	pose_deform_mlp
Expression (τ)	shape_deform_mlp

Table 1. Each primary control vector, and the corresponding MLP module

The loss is then calculated on these captured representations. The combined loss function encourages the representations for gaze and pose to be orthogonal to the representation for expression:

$$L_{ortho_total} = w_1 |\cos_sim(v_{gaze}, v_{expr})| + w_2 |\cos_sim(v_{pose}, v_{expr})| \quad (12)$$

where v_{gaze} , v_{pose} , and v_{expr} are the captured internal representations, and w_1, w_2 are weighting factors. This new

loss term is added to the final training objective. During backpropagation, its gradients update the weights of the first layers of the deformation MLPs, directly training them to learn disentangled projections with negligible computational overhead.

3.6. Training and Loss Functions

The overall training process involves two distinct stages: an initial stage to learn a robust geometric prior, followed by the end-to-end training of the full redirection model.

The **Initialization Stage** follows the methodology of GazeGaussian [24] and Gaussian Head Avatar [26]. An initial optimization is performed to learn a Signed Distance Function (SDF) based neutral geometry, along with face deformation and eye rotation fields from the training data. A neutral mesh is then extracted from this SDF using Deep Marching Tetrahedra (DMTet) [20]. This pre-trained mesh and its associated MLPs provide the initial positions and weights for the Gaussians in our main model, ensuring a stable starting point.

After initializing the Gaussians, the complete redirection pipeline—including the two-stream 3DGS model, the transformation fields, our DiT-based renderer, and the new orthogonality constraint—is trained jointly in an end-to-end fashion. The training is guided by three primary loss functions.

Image Synthesis Loss (L_I): This loss ensures the perceptual quality of the generated images and is a composite of ℓ_1 , SSIM, and LPIPS losses applied to the rendered images of the face-only, eyes, and full head regions. For any given region, the loss is formulated as

$$\mathcal{L}_I^e = \|I_{gt} - I_e\|_1 + \lambda_{SSIM}(1 - SSIM(I_{gt}, I_e)) + \lambda_{VGG} VGG(I_{gt}, I_e) \quad (13)$$

Gaze Redirection Loss (L_G): This functional loss targets redirection accuracy. Following GazeGaussian, we use a pre-trained VGG-based gaze estimator $\psi^g(\cdot)$ to measure the angular error between the gaze estimated from the rendered image I_h and the gaze from the ground-truth target image I_{gt} . For our weak supervision strategy, the ground-truth gaze is replaced with the synthetic intermediate gaze vector. The loss is defined as

$$\mathcal{L}_G(I_h, I_{gt}) = E_{ang}(\psi^g(I_h), \psi^g(I_{gt})) \quad (14)$$

Orthogonality Constraint Loss ($\mathcal{L}_{ortho_total}$): This term encourages disentanglement between the control vectors for gaze, pose, and expression, as previously described.

The **Final Training Objective** is to minimize a weighted sum of these three loss components. The complete loss function for our DiT-Redirection model is

$$\mathcal{L} = \lambda_I \mathcal{L}_I + \lambda_G \mathcal{L}_G + \lambda_{ortho} \mathcal{L}_{ortho_total} \quad (15)$$

where λ_I , λ_G , and λ_{ortho} are weighting hyperparameters.

4. Experiments

4.1. Experimental Settings

Datasets & Preprocessing: The model will be trained on the **ETH-XGaze dataset** [29] and evaluated on its person-specific test set for within-dataset comparison, as well as on the **ColumbiaGaze** [23], **MPIIFaceGaze** [27, 28], and **GazeCapture** [10] datasets for cross-dataset evaluation. We follow GazeNeRF’s [19] standard procedure which is also followed by GazeGaussian [24].

Our preprocessing pipeline begins with the raw images, which undergo **normalization** and are resized to a standard 512×512 resolution. To enable independent rendering of facial features, we generate **segmentation masks** for the face and eye regions using a face parsing model [31]. Concurrently, we apply the 3D face tracking method from [25] to extract the identity codes, expression parameters, and camera poses that serve as inputs to our model. Finally, for consistency across all data sources, we convert the provided gaze labels into **pitch–yaw angles** relative to the head’s coordinate system.

Baselines: Performance will be compared against the **3DGS GazeGaussian** [24] model to measure the impact of our enhancements, alongside other strong baselines such as **GazeNeRF** [19] and **STED** [30].

Evaluation Metrics: To ensure a fair comparison, performance will be measured using the proven comprehensive metrics from GazeGaussian, including **Redirection Accuracy** (Gaze/Head Error), **Image Quality** (SSIM, PSNR, LPIPS, FID), and **Identity Preservation** (ID). **Rendering Speed (FPS)** will be omitted due to a more high-end GPU used compared to baselines, making this comparison unfair.

4.2. Implementation Details

Model Architecture: The Diffusion Transformer (DiT) [16] backbone is implemented with a depth of **6** blocks and **8** attention heads per block. The DiT operates in a latent space, using a pre-trained VAE to encode the input feature map into a **4-channel** latent representation, which is then patchified and fed to the transformer.

Training Hyperparameters: The model is trained using the Adam optimizer [9] with an initial learning rate of **$1e-4$** . The learning rate is managed by a customized step-decay schedule.

Hardware and Software: All models were implemented in PyTorch. Due to the significant computational requirements of the Diffusion Transformer (DiT) architecture, all training and evaluation experiments were conducted on a single **NVIDIA A100-SXM4 GPU**.

4.3. Within-Dataset Comparison

Following the established experimental setup of GazeNeRF and GazeGaussian, we first perform a comprehensive

within-dataset evaluation to benchmark the performance of our proposed method, DiTGaze, against other state-of-the-art models. To ensure a direct and fair comparison, all models are trained on the identical dataset derived from the ETH-XGaze training set, which consists of 14,400 images covering 80 subjects. This dataset is constructed by selecting 10 frames per subject, with each frame providing 18 distinct camera views.

The evaluation is conducted on the person-specific test set of ETH-XGaze, which comprises 15 subjects not seen during training, each with 200 images annotated with precise gaze and head pose labels. We adhere strictly to the pairing protocol defined in GazeNeRF, where these 200 images are paired as input and target samples for the redirection task. This consistent pairing is used across all evaluated models to guarantee a fair assessment.

Table 2 presents the quantitative results of DiTGaze alongside the baseline methods. The results demonstrate that our model achieves a new state-of-the-art, outperforming the GazeGaussian baseline on 6 out of 7 key metrics.

Fidelity and Realism: Our primary contribution, the replacement of the U-Net renderer with a Latent Diffusion Transformer (DiT), yields substantial improvements across all fidelity metrics. We achieve a **1.8 dB gain in PSNR** (20.512 vs. 18.734) and a significant **13.4% reduction in LPIPS** (0.187 vs. 0.216). This is further corroborated by a clear improvement in FID (38.319 vs. 41.972). This demonstrates the DiT’s superior self-attention mechanism for capturing long-range dependencies and synthesizing higher-fidelity, more realistic facial textures than the convolutional baseline, which can also be seen in Figure 3 as our model preserves finer details.

Disentanglement and Accuracy: Our novel Orthogonality Constraint Loss proves highly effective, boosting the Identity Similarity score by 4 points (71.724 vs. 67.749). This confirms our hypothesis that explicitly enforcing the disentanglement of gaze, pose, and shape representations is critical for preserving subject identity. Furthermore, our Intermediate Gaze Sampler successfully refines the model’s understanding of the gaze manifold, reducing the primary **Gaze error to a new SOTA of 6.353°**.

Head Pose: While our model shows a minor regression of 0.22° in Head Pose error compared to GazeGaussian, our score of 2.349° remains highly competitive and is still a significant improvement over prior art such as GazeNeRF (3.470°). This minor regression is a predictable and acceptable trade-off. DiTGaze is trained on a more complex, multi-objective loss function, including a novel weak supervision task from our Intermediate Gaze Sampler. While this new objective successfully improved our gaze accuracy, it creates a more challenging optimization landscape for the 3D Gaussian fields. We conclude that this negligible trade-off in pose stability is massively outweighed by the state-of-

Table 2. Quantitative within-dataset evaluation of our proposed model against state-of-the-art baselines on the ETH-XGaze test set. Metrics include redirection accuracy (gaze and head pose error in degrees), image fidelity (SSIM, PSNR, LPIPS, FID), and Identity Preservation. All trained on ETH-XGaze training set, tested on ETH-XGaze test set.

Method	Gaze↓	Head Pose↓	SSIM↑	PSNR↑	LPIPS↓	FID↓	Identity Similarity↑
STED	16.217	13.153	0.726	17.530	0.300	115.020	24.347
GazeNeRF	6.944	3.470	0.733	15.453	0.291	81.816	45.207
GazeGaussian	6.622	2.128	0.823	18.734	0.216	41.972	67.749
DiT Gaze (Ours)	6.353	2.349	0.850	20.512	0.187	38.319	71.724

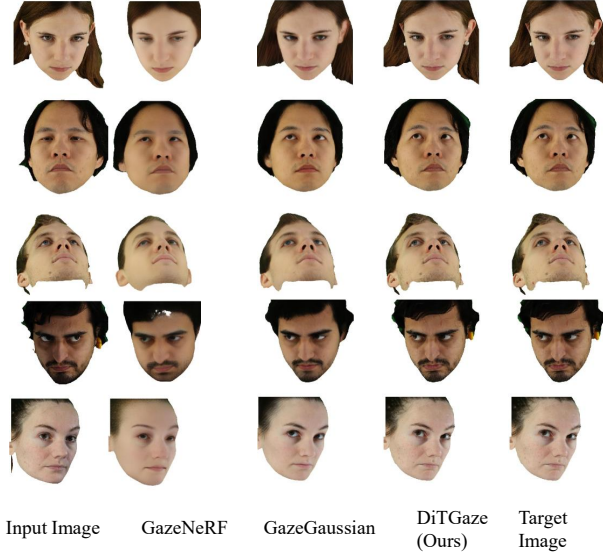


Figure 3. Within-dataset visualization: Head images are generated from the ETH-XGaze test set comparing DiT Gaze (Ours) against GazeNeRF and GazeGaussian. Baseline results are reproduced directly from GazeGaussian’s publication. Our DiT-based model not only preserves identity and matches the target gaze, but also generates superior, high-fidelity facial details, such as individual strands of hair. In contrast, GazeNeRF suffers from significant identity loss and blur, while the GazeGaussian baseline produces softer, less realistic textures.

the-art gains in fidelity, identity preservation, and primary gaze accuracy.

4.4. Cross-Dataset Comparison

To assess the generalization capability of our model, we conduct a rigorous cross-dataset evaluation. All methods are trained on the ETH-XGaze dataset, following the exact same training setup and utilizing the same model parameters as in the within-dataset comparison. The models are then evaluated on three unseen datasets: **ColumbiaGaze**, **MPIIFaceGaze**, and the test set of **GazeCapture**. This protocol is designed to test the models’ robustness and adaptability to variations in subjects, lighting conditions,

and camera setups that are not present in the training data.

The quantitative results, presented in Table 3, confirm that **DiT Gaze robustly generalizes** to unseen data, outperforming the state-of-the-art **GazeGaussian** baseline on the metrics most critical to generalization: fidelity and identity. This superior performance stems directly from our core architectural contributions.

This is most evident in the perceptual quality, where our **DiT renderer** achieves a new state-of-the-art on all three datasets. We observe a significant and consistent reduction in LPIPS error, for example, dropping from 0.273 to **0.231** on **ColumbiaGaze** and from 0.209 to **0.181** on **GazeCapture**. This proves that the DiT’s self-attention mechanism is fundamentally better at synthesizing realistic, high-fidelity textures on novel subjects than the convolutional U-Net baseline.

Furthermore, our results demonstrate the critical impact of our regularization strategies. Our **Orthogonality Constraint Loss** proves highly effective at learning a truly disentangled identity representation. This is validated by a massive **+5.2-point gain** in Identity Similarity on ColumbiaGaze (64.963 vs. 59.788) and consistent, significant gains on both MPIIFaceGaze and GazeCapture. This shows our model is far superior at preserving the identity of unseen subjects. Similarly, our **Intermediate Gaze Sampler** provides a more robust gaze manifold, leading to the lowest Gaze error on all three datasets (e.g., 10.512 vs. 10.943 on MPIIFaceGaze).

However, it is noted that these state-of-the-art gains are balanced by a minor, consistent trade-off in Head Pose stability, which mirrors our findings in the within-dataset evaluation. This suggests an acceptable trade-off, where the complex, multi-objective training (driven by the L_{Ortho} and Gaze Sampler) slightly impacts pose stability in exchange for massive, generalizing improvements in fidelity, identity preservation, and gaze accuracy.

4.5. Ablation Studies

To validate the effectiveness and individual contributions of our proposed components, we conduct a series of ablation experiments on the ETH-XGaze dataset. We systematically dismantle our full model by removing one contribution at a time and retraining the model under the same protocol as

Table 3. Cross-dataset evaluation of our method against state-of-the-art baselines. We report quantitative results on the ColumbiaGaze, MPIIFaceGaze, and GazeCapture datasets, measuring gaze and head pose redirection errors ($^{\circ}$), ID and LPIPS score. All trained on ETH-XGaze training set and tested on the respective datasets.

Method	ColumbiaGaze				MPIIFaceGaze				GazeCapture			
	Gaze↓	Pose↓	LPIPS↓	ID↑	Gaze↓	Pose↓	LPIPS↓	ID↑	Gaze↓	Pose↓	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	6.808
GazeNeRF	9.464	3.811	0.352	23.157	14.933	7.118	0.272	30.981	10.463	9.064	0.232	20.981
GazeGaussian	7.415	3.332	0.273	59.788	10.943	5.685	0.224	41.505	9.752	7.061	0.209	19.025
DiTGaze (Ours)	7.265	3.611	0.231	64.963	10.512	5.905	0.196	44.478	9.676	7.217	0.181	22.236

Table 4. Component Wise ablation study of our proposed components on the ETH-XGaze dataset. We analyze the impact of the DiT renderer, orthogonality loss, and weak supervision strategy on redirection errors (gaze, head pose) and image quality (LPIPS, FID).

DiT	Orthogonality Loss	Intermediate Gaze Weak Supervision	Gaze↓	Pose↓	LPIPS↓	FID↓
	✓	✓	6.485	2.328	0.214	41.878
✓		✓	6.396	2.572	0.191	39.337
✓	✓		6.714	2.344	0.190	38.582
✓	✓	✓	6.353	2.349	0.187	38.319

before. The results are presented quantitatively in Table 4.

w/o DiT Renderer. In this variant, we replace our DiT Renderer with the baseline convolutional U-Net from GazeGaussian. The results, shown in the first row, are unambiguous: the fidelity metrics worsen. LPIPS error increases by **14.4%** ($0.187 \rightarrow 0.214$) and FID degrades by **9.3%** ($38.319 \rightarrow 41.878$), reverting both metrics close to the original GazeGaussian baseline scores. This directly validates that our DiT’s self-attention mechanism is the primary driver of our model’s state-of-the-art fidelity.

w/o Weak Supervision. To verify the benefit of our Intermediate Gaze Sampler (third row), we trained a version of our model using only the dataset’s discrete source-target pairs. The impact is significant: the Gaze error increases by **5.7%** ($6.353 \rightarrow 6.714$). This result is notably worse than even the GazeGaussian baseline (6.622), proving that our weak supervision strategy is essential for learning a smooth, continuous gaze manifold and achieving our final state-of-the-art accuracy.

w/o Orthogonality Loss. When we remove our Orthogonality Constraint Loss (second row), the Head Pose error increases by **9.5%** ($2.349 \rightarrow 2.572$), becoming the worst of any ablation. This confirms our hypothesis that explicitly enforcing disentanglement within the Face Deform Field is critical for maintaining geometric stability and preventing pose from being corrupted during the redirection task.

In summary, the ablation experiments clearly demonstrate that each of our three contributions provides distinct benefits. The full model’s superior performance is a direct and measurable result of the **DiT’s** advanced synthesis capabilities, the **Orthogonality Loss’s** enhancement of pose stability, and the **Gaze Sampler’s** improvement to gaze accuracy.

5. Conclusion and Discussion

We presented **DiTGaze**, a framework that advances Gaze Redirection with three novel contributions: (1) a **Latent Diffusion Transformer (DiT)**, (2) an **Orthogonality Constraint Loss** for superior disentanglement, and (3) an **Intermediate Gaze Sampler** to learn a smooth, continuous gaze manifold.

Our comprehensive experiments conclusively demonstrate the effectiveness of our contributions, establishing a new state-of-the-art. On ETH-XGaze, DiTGaze outperforms GazeGaussian [1] with a **1.8 dB PSNR gain** and a **13.4% LPIPS reduction**, which we attribute to the DiT’s **self-attention** mechanism for global coherence and **AdaLN** conditioning for subject-specific detail. This superiority extends to generalization, where our model shows significant and consistent gains in **Identity Similarity** and **LPIPS** across all three cross-dataset benchmarks. Furthermore, our ablation study confirms that each contribution is essential, as removing the DiT, Orthogonality Loss, or Gaze Sampler caused a significant degradation in fidelity, pose stability, and gaze accuracy, respectively.

Limitations. DiTGaze provides a superior method for offline synthetic data generation. Its primary limitation is the increased computational cost of the DiT, which slows inference speed. We also note a minor, 0.22° regression in head pose stability, a likely trade-off for our more complex, multi-objective training.

Future Work. Future work could explore model distillation and quantization-aware training to create a lightweight, real-time version of DiTGaze. Further investigation into the trade-off between gaze accuracy and pose stability could also yield new, more robust regularization strategies.

References

- [1] Timo Bolkart, Tianye Li, Michael J. Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4d scans. *ACM Transactions on Graphics*, 36(6):1–17, 2017. 2
- [2] Y. et al. Choi. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *CVPR*, 2017. 2
- [3] Yaroslav Ganin, Daniil Kononenko, Diana Sungatullina, and Victor Lempitsky. Deepwarp: Photorealistic image resynthesis for gaze manipulation. In *ECCV*, pages 311–326. Springer, 2016. 1, 2
- [4] Zhe He, Adrian Spurr, Xucong Zhang, and Otmar Hilliges. Photo-realistic monocular gaze redirection using generative adversarial networks. In *ICCV*, pages 6931–6940, 2019. 1, 2
- [5] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 2020. 2
- [6] Yuelang Hong, Bing Peng, Hao Xiao, Lingjie Liu, and Juyong Zhang. Headnerf: A realtime nerf-based parametric head model. In *CVPR*, pages 20342–20352, 2021. 2
- [7] S. et al. Jin. Redirtrans: Latent-to-latent translation for gaze and head redirection. In *CVPR*, 2023. 2
- [8] Bernd Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023. 1, 2
- [9] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2015. 6
- [10] Kyle Krafka et al. Eye tracking for everyone. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2176–2184, 2016. 6
- [11] 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. *Communications of the ACM*, 65(1):99–106, 2022. 1
- [12] S. et al. Mo. Dit-3d: Exploring plain diffusion transformers for 3d shape generation. *NeurIPS*, 2023. 2
- [13] T. et al. Nguyen. Climax: A foundation model for weather and climate. In *ICML*, 2023. 2
- [14] Shaoyan Pan, Tonghe Wang, Richard L J Qiu, Marian Axente, Chih-Wei Chang, Junbo Peng, Ashish B Patel, Joseph Shelton, Sagar A Patel, Justin Roper, and Xiaofeng Yang. 2D medical image synthesis using transformer-based denoising diffusion probabilistic model. *Phys. Med. Biol.*, 68(10), 2023. 2
- [15] Sangjin Park, Daeha Kim, and Byung Cheol Song. Fine gaze redirection learning with gaze hardness-aware transformation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3464–3473, 2023. 4
- [16] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *ICCV*, pages 4172–4182, 2023. 2, 3, 6
- [17] Robin et al. Rombach. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2021. 2, 4
- [18] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 2
- [19] Andrea Ruzzi et al. Gazenerf: 3d-aware gaze redirection with neural radiance fields. In *CVPR*, pages 9676–9685, 2022. 1, 2, 6
- [20] Tiancheng Shen, Jun Gao, Kangxue Yin, Ming-Yu Liu, and Sanja Fidler. Deep marching tetrahedra: A hybrid representation for high-resolution 3d shape synthesis. In *NeurIPS*, pages 6087–6101, 2021. 5
- [21] Tianchang Shen, Jun Gao, Kangxue Yin, Ming-Yu Liu, and Sanja Fidler. Deep marching tetrahedra: a hybrid representation for high-resolution 3d shape synthesis. *Advances in Neural Information Processing Systems*, 34:6087–6101, 2021. 3
- [22] Z. Shu, E. Shechtman, D. Samaras, and S. Hadap. Eye-opener: Editing eyes in the wild. *ACM TOG*, 2017. 2
- [23] Brian A Smith, Qin Yin, Steven K Feiner, and Shree K Nayar. Gaze locking: Passive eye contact detection for human-object interaction. In *Proceedings of the 26th Annual ACM Symposium on User Interface Software and Technology*, 2013. 6
- [24] Xiaoyu Wei, Peng Chen, Guanyu Li, Ming Lu, Hui Chen, and Feng Tian. Gazegaussian: High-fidelity gaze redirection with 3d gaussian splatting. *arXiv preprint arXiv:2501.XXXX*, 2025. 1, 2, 3, 5, 6
- [25] Yuelang Xu, Benwang Chen, Zhe Li, Hongwen Zhang, Lizhen Wang, Zerong Zheng, and Yebin Liu. Gaussian head avatar: Ultra high-fidelity head avatar via dynamic gaussians. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1931–1941, 2024. 3, 6
- [26] Yufeng Xu et al. Gaussian head avatar: Ultra high-fidelity head avatar via dynamic gaussians. In *CVPR*, pages 1931–1941, 2024. 2, 3, 5
- [27] Xucong Zhang, Yusuke Sugano, Mario Fritz, and Andreas Bulling. Appearance-based gaze estimation in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4511–4520, 2015. 6
- [28] Xucong Zhang, Yusuke Sugano, Mario Fritz, and Andreas Bulling. It’s written all over your face: Full-face appearance-based gaze estimation. In *CVPR Workshops*, pages 2299–2308, 2016. 6
- [29] Xucong Zhang, Seonwook Park, Thabo Beeler, Derek Bradley, Siyu Tang, and Otmar Hilliges. Eth-xgaze: A large scale dataset for gaze estimation under extreme head pose and gaze variation. In *European Conference on Computer Vision*, 2020. 6
- [30] Yufei Zheng, Seung-Hwan Park, Xucong Zhang, Shalini De Mello, and Otmar Hilliges. Self-learning transformations for improving gaze and head redirection. *NeurIPS*, 33, 2020. 1, 2, 6
- [31] zll. face-parsing.pytorch: Using modified bisenet for face parsing in pytorch. <https://github.com/zllrunning/face-parsing.PyTorch>. Accessed: 2025-11-12. 6