

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

Supplementary Material

805 A. Overview

806 The supplementary material encompasses the subse-
807 quent components. Please visit the anonymous website
808 <https://gazegaussian.github.io/> for additional visual compar-
809 isons of novel view and novel gaze synthesis.

- 810 • Supplementary experiments
 - 811 – Ablation study on cross-dataset
 - 812 – Personal calibration for gaze estimation
 - 813 – Comparison between GazeNeRF + expression-guided
 - 814 neural renderer and GazeGaussian
- 815 • Additional visualization results
 - 816 – Visualization for ablation study
 - 817 – Visualization for cross-dataset comparison
 - 818 – Visualization for identity morphing
 - 819 – Visualization for transformed Gaussians
- 820 • Implementation details
- 821 • Dataset and pre-processing details
- 822 • Ethical consideration and limitations

823 B. Supplementary experiments

824 B.1. Ablation study on cross-dataset

825 To further validate the effectiveness of each proposed com-
826 ponent, we conduct an ablation study on the cross-dataset
827 evaluation to assess the generalization capability of our full
828 pipeline. As shown in Tab. 1, the results are consistent with
829 the ablation study in the main text. The proposed Gaussian
830 eye rotation representation significantly improves eye redi-
831 rection accuracy while ensuring robust redirection across
832 cross-domain datasets. Additionally, the expression-guided
833 neural renderer enhances the fidelity of the synthesized im-
834 ages, preserving the identity characteristics of the input im-
835 age. From the ablation study on cross-dataset, we can further
836 validate the importance of each component.

837 B.2. Personal calibration for gaze estimation

838 Following GazeNeRF, we perform personal calibration to
839 demonstrate the benefits of our method for downstream gaze
840 estimation tasks. Specifically, given a few calibration sam-
841 ples from person-specific test sets, we augment these real
842 samples with gaze-redirectioned samples generated by Gaze-
843 Gaussian. We then fine-tune a gaze estimator pre-trained on
844 ETH-XGaze using these augmented samples and compare
845 its performance with a baseline model fine-tuned only on
846 real samples. To ensure a fair comparison, the total num-
847 ber of augmented samples is fixed at 200 (real + generated
848 samples), and we vary the number of real samples used for
849 fine-tuning during the evaluation.

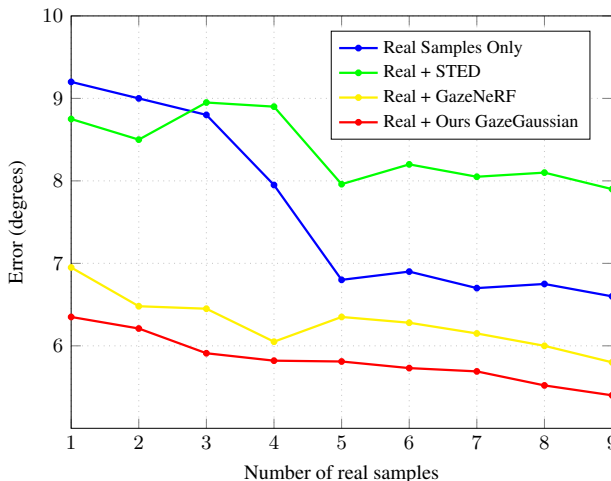


Figure 1. Error comparison based on number of real samples.

850 As shown in Fig. 1, the x-axis represents the number of
851 real samples used, and the y-axis shows the gaze estimation
852 error in degrees on the ETH-XGaze person-specific test set.
853 We evaluate up to nine real samples in the few-shot setting.
854 Fine-tuning the pre-trained gaze estimator with real and gen-
855 erated samples from GazeGaussian achieves the lowest gaze
856 estimation error across all settings. Compared to GazeNeRF,
857 GazeGaussian demonstrates a clear advantage, especially
858 when fewer real samples are available, indicating that the
859 generated samples from GazeGaussian are of higher fidelity
860 and more effective for improving downstream gaze estima-
861 tion accuracy. In contrast, samples generated by GazeNeRF
862 lead to higher errors, while STED performs the worst, show-
863 ing a notable limitation in leveraging 2D generative models
864 for this task. This is due to the lack of consideration for the
865 3D nature of gaze redirection in STED, which is critical for
866 high-quality sample generation and effective downstream
867 adaptation.

868 B.3. Comparison between GazeNeRF + expression- 869 guided neural renderer and GazeGaussian

870 We compare the performance of GazeNeRF, GazeNeRF en-
871 hanced with the expression-guided neural renderer (EGNR),
872 and our proposed GazeGaussian on the ETH-XGaze dataset.
873 As shown in Tab. 2, integrating EGNR into GazeNeRF leads
874 to noticeable improvements in gaze redirection accuracy
875 and image quality. This demonstrates the versatility of the
876 proposed expression-guided neural renderer in enhancing
877 facial synthesis and better capturing identity-specific expres-
878 sions. However, even with the added EGNR, GazeNeRF’s
879 performance remains limited compared to GazeGaussian.

Table 1. Component-wise ablation study of GazeGaussian on the ColumbiaGaze, MPIIFaceGaze and GazeCapture datasets.

Two-stream	Gaus. Eye Rep.	Exp. Guided	ColumbiaGaze				MPIIFaceGaze				GazeCapture			
			Gaze↓	Head↓	LPIPS↓	ID↑	Gaze↓	Head↓	LPIPS↓	ID↑	Gaze↓	Head↓	LPIPS↓	ID↑
✓			8.996	4.494	0.325	49.286	19.787	8.491	0.321	34.483	15.697	13.740	0.260	33.393
✓		✓	9.143	4.509	0.324	49.805	16.689	8.578	0.303	34.194	15.926	14.869	0.261	33.004
✓	✓		7.799	3.754	0.284	57.252	11.938	6.860	0.257	35.614	10.339	8.208	0.216	40.458
✓	✓	✓	7.710	3.899	0.280	58.969	12.559	6.188	0.246	37.444	11.296	8.460	0.224	42.294
✓	✓	✓	7.415	3.332	0.273	59.788	10.943	5.685	0.224	41.505	9.752	7.061	0.209	44.007

Table 2. Comparison between GazeNeRF + expression-guided neural renderer and GazeGaussian on ETH-xgaze

Method	Gaze↓	Head Pose↓	SSIM↑	PSNR↑	LPIPS↓	FID↓	Identity Similarity↑	FPS↑
GazeNeRF	6.944	3.470	0.733	15.453	0.291	81.816	45.207	46
GazeNeRF + EGNR	6.854	3.025	0.764	16.147	0.258	67.219	50.268	44
GazeGaussian (Ours)	6.622	2.128	0.823	18.734	0.216	41.972	67.749	74

880 The fundamental constraint lies in GazeNeRF’s representa- 915
 881 tion, which lacks the explicit modeling of gaze and facial 916
 882 expression dynamics offered by GazeGaussian’s two-stream 917
 883 Gaussian structure. GazeNeRF struggles to achieve fine- 918
 884 grained expression synthesis and accurate gaze alignment, 919
 885 which are critical for high-fidelity gaze redirection. 920

886 In contrast, GazeGaussian leverages the strengths of 921
 887 the expression-guided neural renderer with its specialized 922
 888 Gaussian-based eye rotation representation and two-stream 923
 889 structure, enabling superior expression modeling and gaze 924
 890 control. This allows GazeGaussian to achieve higher fidelity, 925
 891 identity preservation, and rendering accuracy compared to 926
 892 GazeNeRF, even when enhanced with the expression-guided 927
 893 neural renderer. These results highlight the importance of 928
 894 combining advanced neural rendering techniques with a ro- 929
 895 bust facial and eye modeling framework for state-of-the-art 930
 896 performance. 931

897 C. Supplementary visualization 932

898 C.1. Visualization for ablation study 933

899 Fig. 2 presents additional qualitative results from our ab- 934
 900 lation study conducted on the ETH-XGaze dataset. These 935
 901 visualizations highlight the importance of each proposed 936
 902 component in GazeGaussian. 937

903 Without the Gaussian eye rotation representation, the 938
 904 model struggles to achieve accurate eye control, resulting in 939
 905 noticeable deviations in gaze direction and reduced realism 940
 906 in the eye region. This demonstrates the critical role of the 941
 907 Gaussian eye rotation representation in enabling precise and 942
 908 realistic gaze redirection. Additionally, the absence of the 943
 909 expression-guided neural renderer leads to a significant loss 944
 910 in facial detail and expression fidelity. With the renderer in- 945
 911 cluded, the synthesized images exhibit finer facial details and 946
 912 improved consistency with the target identity, showcasing 947
 913 the renderer’s effectiveness in enhancing the overall quality 948
 914 of face synthesis. These results confirm that both compo-

nents contribute significantly to the superior performance 915
 and visual fidelity of GazeGaussian. 916

917 C.2. Visualization for cross-dataset comparison 918

919 We provide additional cross-dataset comparison visualiza- 920
 921 tions for MPIIFaceGaze (Fig. 4), ColumbiaGaze (Fig. 5) and 922
 923 GazeCapture (Fig. 6). Compared to the baseline, GazeGaus- 924
 925 sian achieves high-fidelity gaze redirection with superior 926
 927 image synthesis quality. 928

929 C.3. Visualization for identity morphing 930

931 Fig. 7 showcases identity morphing results on the ETH- 932
 933 XGaze dataset. For this experiment, we randomly select 934
 935 two subjects with identical gaze directions and head poses. 936
 937 By interpolating their latent codes, we generate a smooth 938
 939 transition between the two identities while keeping the gaze 940
 941 direction and head pose consistent. 942

943 This visualization demonstrates the capability of Gaze- 930
 944 Gaussian to preserve gaze alignment and head orientation 931
 945 during synthesis, even as the facial features gradually change 932
 946 according to the interpolated latent codes. The results high- 933
 947 light the robustness of GazeGaussian in maintaining high- 934
 948 fidelity gaze redirection while adapting facial characteristics 935
 949 as required. This ability to control identity-specific details 936
 950 while preserving gaze and pose consistency underscores the 937
 951 flexibility and effectiveness of the proposed method. 938

939 C.4. Visualization for transformed Gaussians 940

941 To demonstrate the advantages of GazeGaussian’s explicit 942
 943 incorporation of head pose and gaze direction for rotating 943
 944 Gaussians in the head and eye regions, we visualize the Gaus- 944
 945 sians after deformation from the canonical space. As shown 945
 946 in Fig. 3, the explicit support for rotation and translation 946
 947 in GazeGaussian allows the deformed Gaussians to form a 947
 948 reasonable spatial distribution and accurate color representa- 948
 949 tion. This capability enables precise geometric control and 949

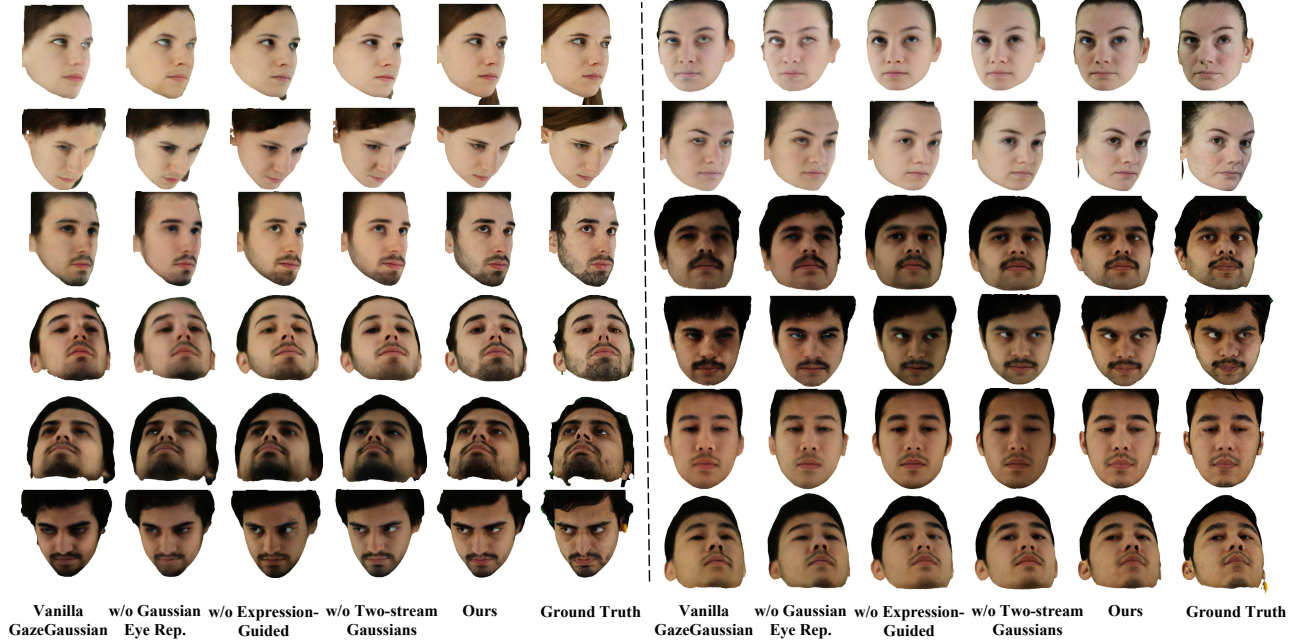


Figure 2. Additional qualitative ablation study on the ETH-XGaze dataset.

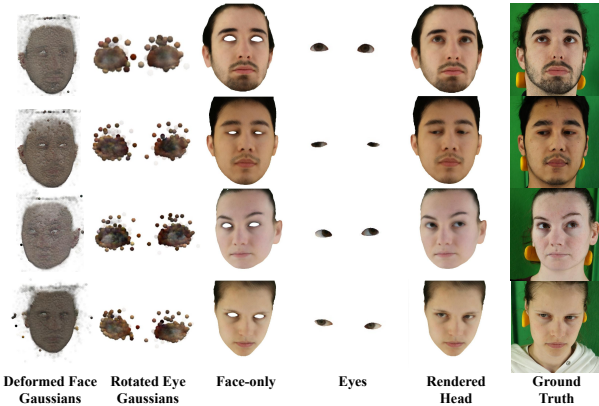


Figure 3. Visualization of transformed two-stream Gaussians after deformation from the canonical space.

948 high-fidelity image rendering. In contrast, GazeNeRF per-
 949 forms rotations only on the feature map level, failing to fully
 950 deform in 3D space, which limits its performance compared
 951 to our method.

952 D. Implementation details

953 We use the Adam optimizer [19], with a learning rate that
 954 follows an exponential decay schedule, starting at 1×10^{-4} .
 955 We use the VGG-based network pre-trained on ImageNet, as
 956 provided by the GazeNeRF [36] implementation, and fine-
 957 tune it on the ETH-XGaze training set for the functional
 958 loss \mathcal{L}_G as the pre-trained gaze estimator. Additionally, we
 959 utilize the ResNet50 backbone from the GazeNeRF [36]
 960 framework, trained on the ETH-XGaze training set, to output
 961 gaze and head pose for evaluation purposes. All experiments

were conducted on an NVIDIA 4090 GPU. We first train
 an SDF network to extract the neutral mesh and initialize
 the two-stream Gaussian parameters in 10 epochs. The full
 pipeline was then trained for an additional 20 epochs until
 convergence. The loss weights follow the same configuration
 as described in the method section of the main text.

E. Dataset and pre-processing details

Following the baseline GazeNeRF [36], all experiments are
 conducted on four widely used datasets.

ETH-XGaze [59] is a large-scale gaze estimation dataset
 featuring high-resolution images across a wide range of
 head poses and gaze directions. Captured with a multi-view
 camera setup under varying lighting conditions, it includes
 756,000 frames from 80 subjects for training. Each frame
 contains images from 18 different camera perspectives. Ad-
 ditionally, a person-specific test set includes 15 subjects,
 each with 200 images provided with ground-truth gaze data.
ColumbiaGaze [39] contains 5,880 high-resolution images
 from 56 subjects. For each subject, images were taken in
 five distinct head poses, with each pose covering 21 preset
 gaze directions, allowing for detailed gaze estimation in con-
 trolled conditions.

MPIIFaceGaze [56, 57] is tailored for appearance-based
 gaze prediction. MPIIFaceGaze offers 3,000 face images for
 each of 15 subjects, paired with two-dimensional gaze labels
 to facilitate gaze estimation research.

GazeCapture [21] is a large-scale dataset collected through
 crowd-sourcing, featuring images captured across different
 poses and lighting conditions. For cross-dataset comparison,
 we use only the test portion, which includes data from 150

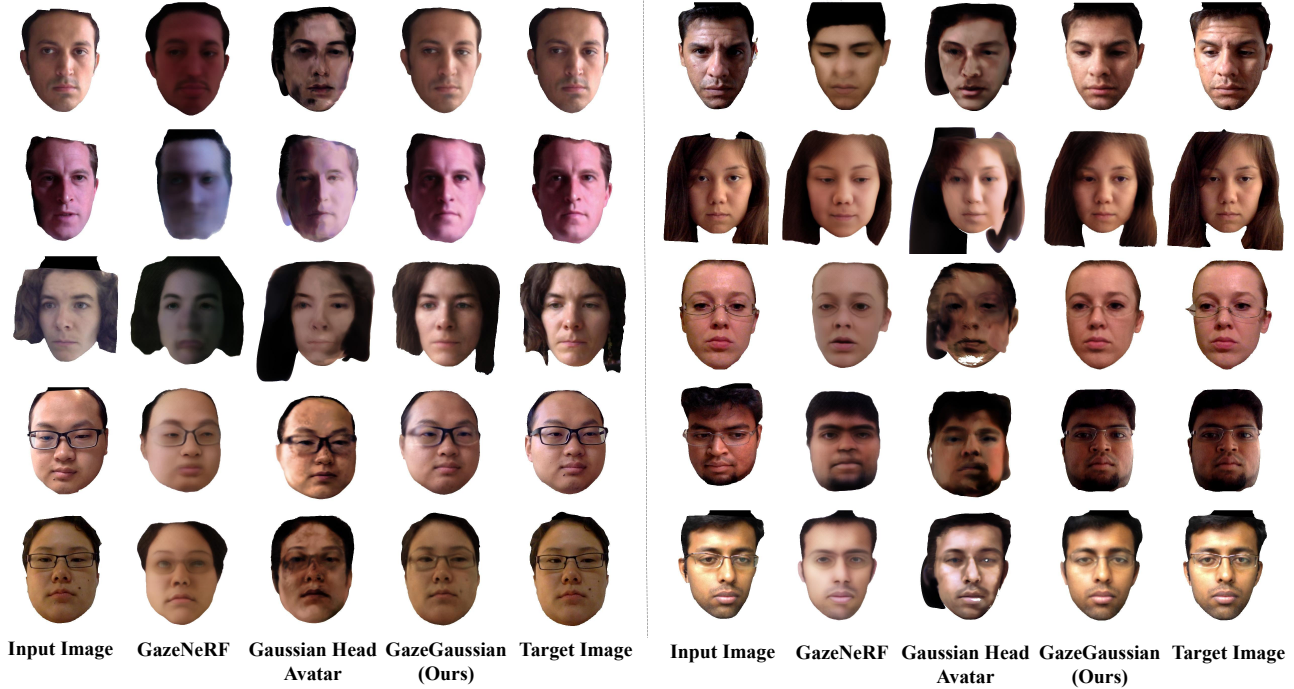


Figure 4. Cross-dataset comparison: Visualization of generated images from the MPIIFaceGaze using our GazeGaussian, GazeNeRF, and Gaussian Head Avatar.

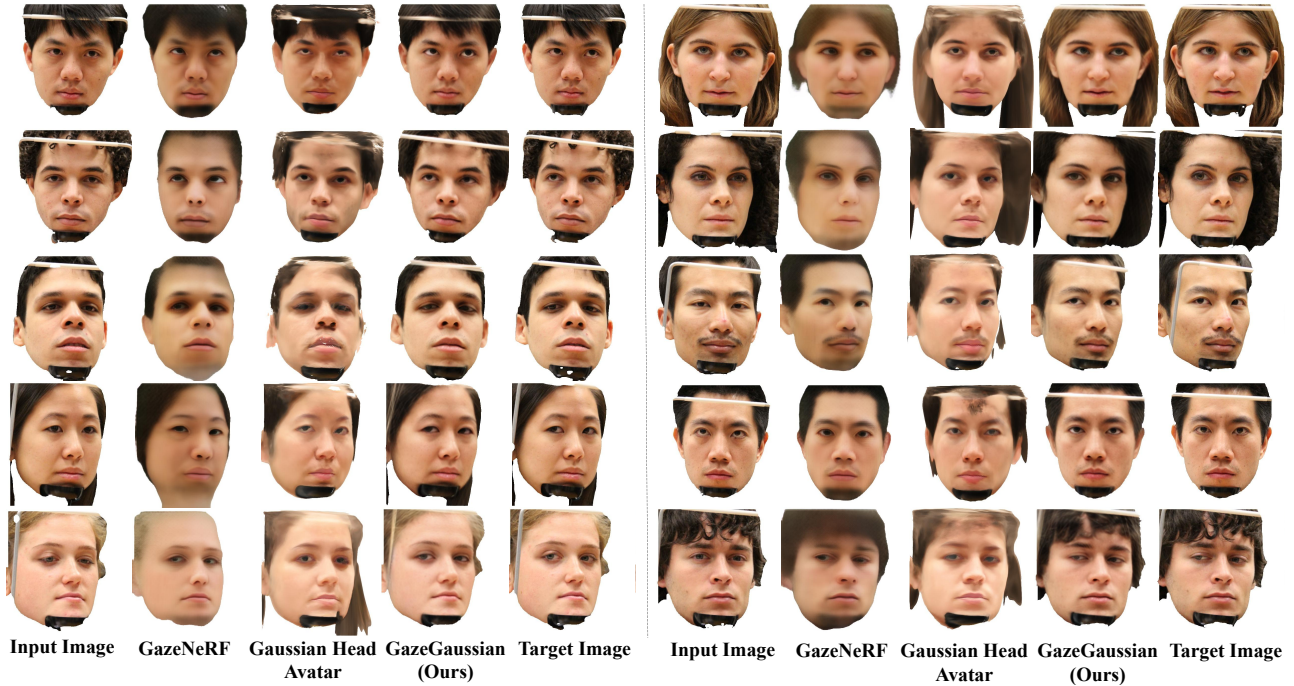


Figure 5. Cross-dataset comparison: Visualization of generated images from the ColumbiaGaze using our GazeGaussian, GazeNeRF, and Gaussian Head Avatar.

992 distinct subjects.

993 **Pre-processing.** We follow the preprocessing steps in
994 GazeNeRF [36] and Gaussian Head Avatar [51]. The origi-
995 nal resolution of ETH-XGaze [59] images is $6K \times 4K$, while
996 images from other datasets vary in resolution. To standardize,

we preprocess all images using the normalization method,
aligning the rotation and translation between the camera and
face coordinate systems. The normalized distance from the
camera to the face center is fixed at 680mm. To extract
3DMM parameters and generate masks for the eyes and face-

997
998
999
1000
1001

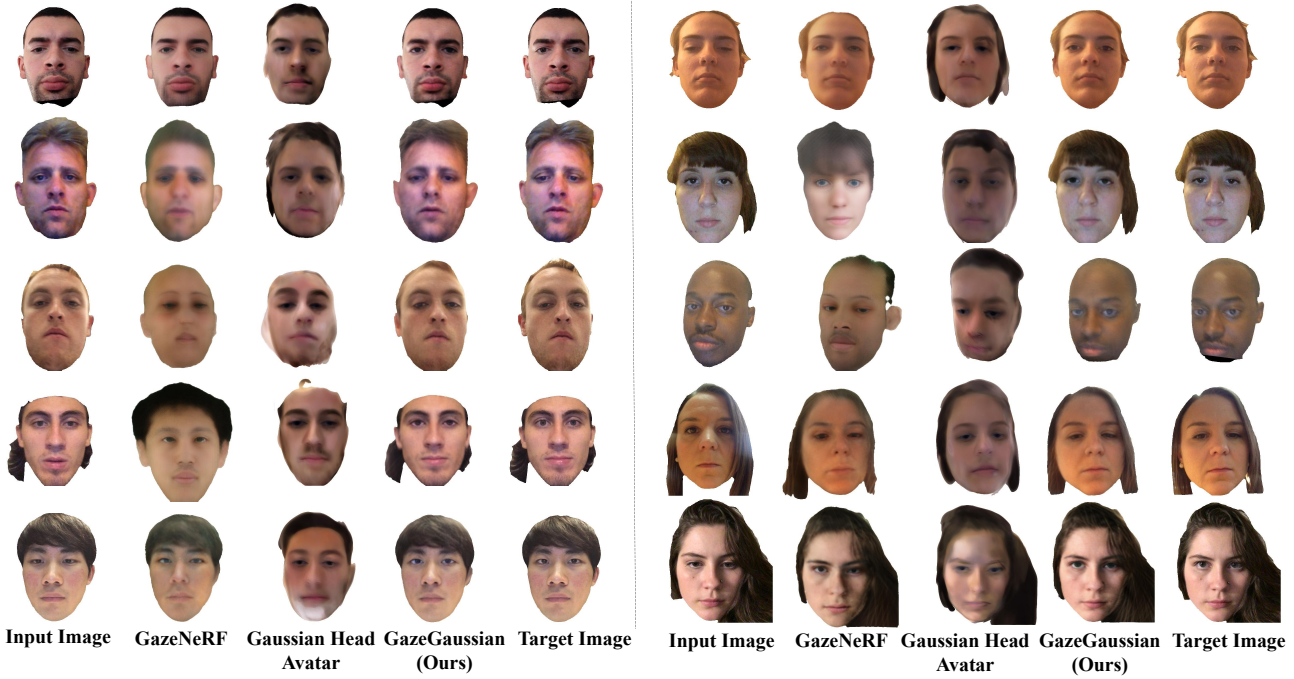


Figure 6. Cross-dataset comparison: Visualization of generated images from the GazeCapture using our GazeGaussian, GazeNeRF, and Gaussian Head Avatar.

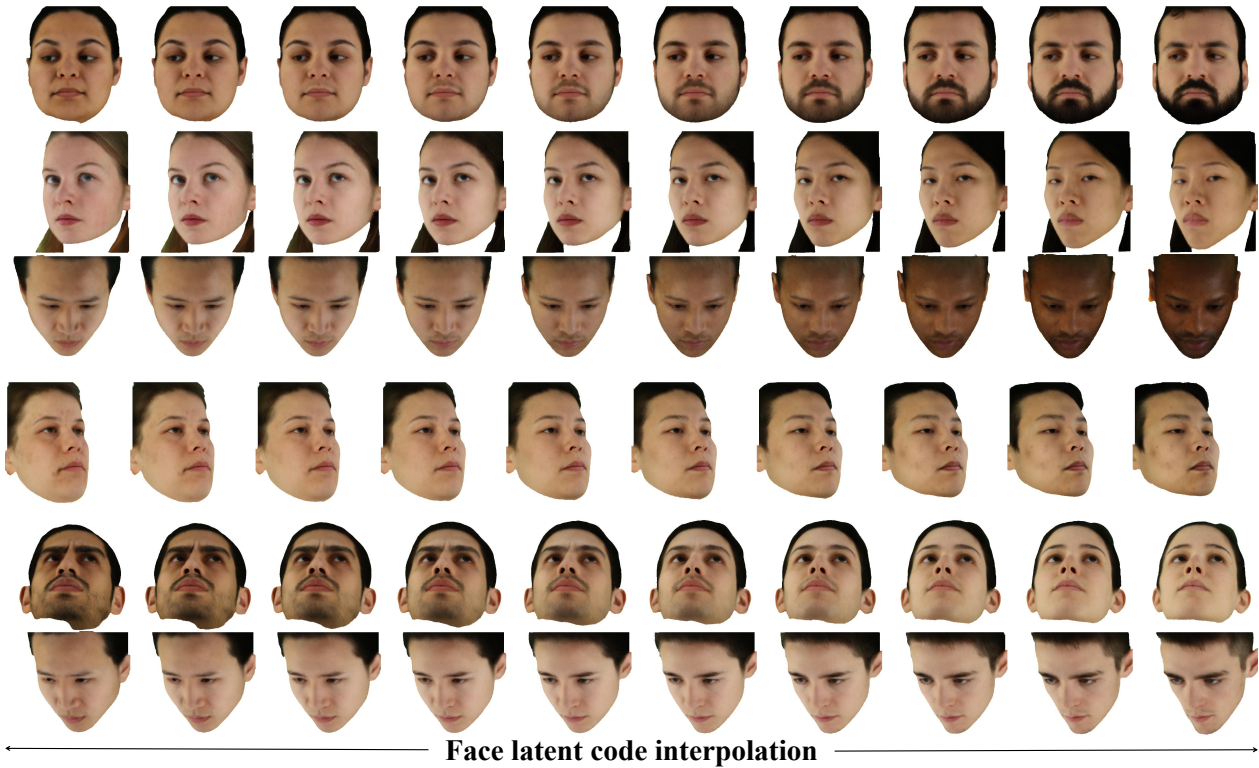


Figure 7. Face morphing results on the ETH-XGaze dataset.

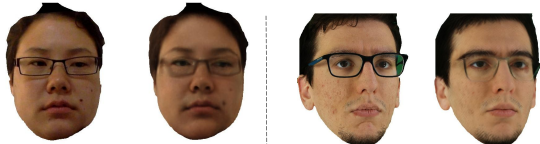
1002 only regions, we utilize the face parsing model from [64].
 1003 GazeGaussian is trained on a single NVIDIA 4090 GPU
 1004 for 20 epochs on the train set from ETH-XGaze. During
 1005 inference, GazeGaussian fine-tunes on a single input image,

taking approximately 30 seconds for fine-tuning and 0.2
 seconds per image for generation.

1006
 1007

1008 F. Ethical consideration and limitations

1009 Our method enables the generation of highly realistic portrait
1010 videos, which, if misused, could contribute to the spread of
1011 misinformation, manipulate public opinion, and undermine
1012 trust in media sources, with significant societal consequences.
1013 Therefore, it is essential to develop reliable methods to dif-
1014 ferentiate between authentic and fabricated content. We
1015 strongly condemn the unauthorized or malicious use of this
1016 technology and emphasize the importance of considering
1017 ethical implications in its deployment.



Target Image GazeGaussian Target Image GazeGaussian

Figure 8. Example of a failure case.

1018 While GazeGaussian represents a significant advance-
1019 ment in gaze redirection quality, there is still one unresolved
1020 issue. Due to limitations in facial tracking models such as
1021 FLAME, it remains challenging to accurately model acces-
1022 sories such as glasses, earrings, and even hair details as
1023 shown in Fig. 8. An existing method [26] has attempted
1024 to use cylindrical Gaussian representations to capture the
1025 movement of long hair. To further enhance the diversity of
1026 character generation, improving the 3DGS facial representa-
1027 tion will be a key focus of our future work.