

3D-aware Blending with Generative NeRFs

Hyunsu Kim¹Gayoung Lee¹Yunjey Choi¹Jin-Hwa Kim^{1,2}Jun-Yan Zhu³¹NAVER AI Lab²SNU AIIS³CMU

Abstract

Image blending aims to combine multiple images seamlessly. It remains challenging for existing 2D-based methods, especially when input images are misaligned due to differences in 3D camera poses and object shapes. To tackle these issues, we propose a 3D-aware blending method using generative Neural Radiance Fields (NeRF), including two key components: 3D-aware alignment and 3D-aware blending. For 3D-aware alignment, we first estimate the camera pose of the reference image with respect to generative NeRFs and then perform pose alignment for objects. To further leverage 3D information of the generative NeRF, we propose 3D-aware blending that utilizes volume density and blends on the NeRF’s latent space, rather than raw pixel space. Collectively, our method outperforms existing 2D baselines, as validated by extensive quantitative and qualitative evaluations with FFHQ and AFHQ-Cat. Please find the code and data on our project [page](#).

1. Introduction

Image blending aims at combining elements from multiple images naturally, enabling a wide range of applications in content creation, and virtual and augmented realities [109, 110]. However, blending images seamlessly requires delicate adjustment of color, texture, and shape, often requiring users’ expertise and tedious manual processes. To reduce human efforts, researchers have proposed various automatic image blending algorithms, including classic methods [72, 58, 10, 87] and deep neural networks [107, 91, 64].

Despite significant progress, blending two unaligned images remains a challenge. Current 2D-based methods often assume that object shapes and camera poses have been accurately aligned. As shown in Figure 1c, even slight misalignment can produce unnatural results, as it is obvious to human eyes that foreground and background objects were captured using different cameras. Several methods [39, 61, 15, 77, 99] warp an image via 2D affine transformation. However, these approaches do not account for 3D geometric differences, such as out-of-plane rotation and 3D shape differences. 3D alignment is much more difficult for users and algorithms,

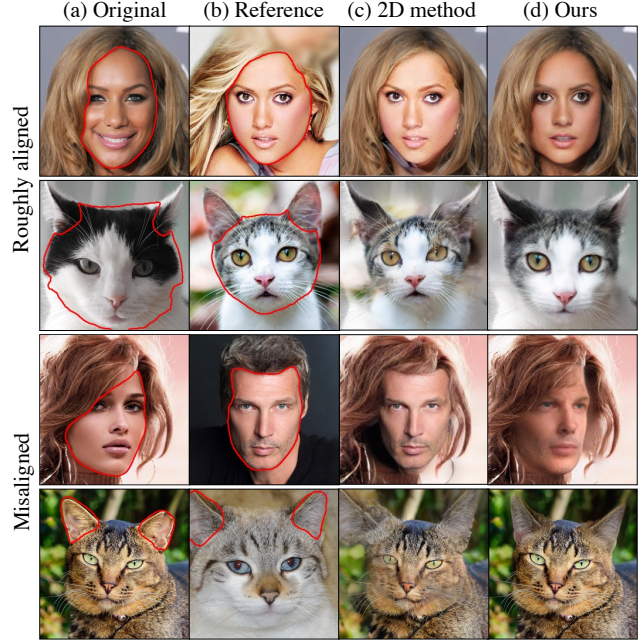


Figure 1: Image blending is challenging for unaligned original and reference images. Existing 2D-based methods [50] struggle to synthesize realistic results due to the 3D object pose differences between foreground and background. In contrast, we propose a 3D-aware blending method that aligns and composes unaligned images without manual effort.

as it requires inferring the 3D structure from a single view. Additionally, even though previous methods get aligned images, they blend images in 2D space. Blending images using only 2D signals, such as pixel values (RGB) or 2D feature maps, doesn’t account for the 3D structure of objects.

To address the above issues, we propose a 3D-aware image blending method based on generative Neural Radiance Fields (NeRFs) [12, 37, 13, 69, 78, 105]. Generative NeRFs learn to synthesize images in 3D using only collections of single-view images. Our method projects the input images to the latent space of generative NeRFs and performs 3D-aware alignment by novel view synthesis. We then perform blending on NeRFs’ latent space. Concretely, we formulate an optimization problem in which a latent code is optimized

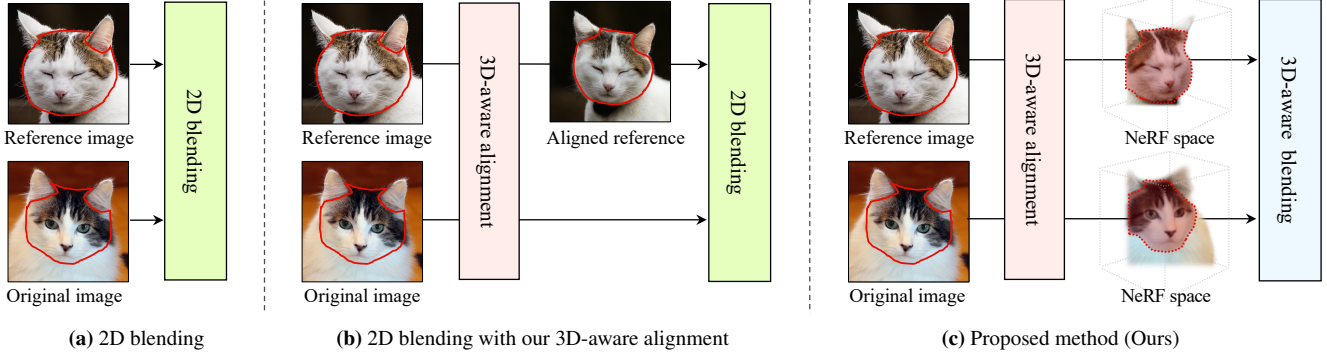


Figure 2: **Comparison with the existing blending methods.** Red lines denote target blending parts. **(a) 2D blending.** 2D blending methods compose two images without any 3D-aware alignment. **(b) 2D blending with 3D-aware alignment.** To address misalignment, we apply our 3D-aware alignment method to existing 2D blending methods. **(c) Proposed method.** We propose 3D-aware blending after applying our 3D-aware alignment. Note that all methods do not use 3D labels or 3D morphable models.

to synthesize an image and volume density of the foreground close to the reference while preserving the background of the original.

Figure 2 shows critical differences between our approach and previous methods. Figure 2a shows a classic 2D blending method composing two 2D images without alignment. We then show the performance of the 2D blending method can be improved using our 3D-aware alignment with generative NeRFs as shown in Figure 2b. To further exploit 3D information, we propose to compose two images in the NeRFs’ latent space instead of 2D pixel space. Figure 2c shows our final method.

We demonstrate the effectiveness of our 3D-aware alignment and 3D-aware blending (volume density) on unaligned images. Extensive experiments show that our method outperforms both classic and learning-based methods regarding both photorealism and faithfulness to the input images. Additionally, our method can disentangle color and geometric changes during blending, and create multi-view consistent results. To our knowledge, our method is the first general-purpose 3D-aware image blending method capable of blending a diverse set of unaligned images.

2. Related Work

Image blending aims to compose different visual elements into a single image. Seminal works tackle this problem using various low-level visual cues, such as image gradients [72, 42, 86, 32, 85], frequency bands [10, 9], color and noise transfer [94, 83], and segmentation [58, 75, 2, 60]. Later, researchers developed data-driven systems to compose objects with similar lighting conditions, camera poses, and scene contexts [59, 18, 39].

Recently, various learning-based methods have been proposed, including blending deep features instead of pix-

els [84, 24, 35] or designing loss functions based on deep features [101, 100]. Generative Adversarial Networks (GAN) have also been used for image blending [91, 106, 25, 50, 11, 109, 79]. For example, In-DomainGAN [106] exploits GAN inversion to achieve seamless blending, and StyleMapGAN [50] blends images in the spatial latent space. Recently, SDEdit [64] proposes a blending method via diffusion models. The above learning-based methods tend to be more robust than pixel-based methods. But given two images with large pose differences, both may struggle to preserve identity or generate unnatural effects.

In specific domains like faces [95, 26, 66, 93] or hair [109, 110, 23, 52], multiple methods can swap and blend unaligned images. However, these methods are limited to faces or hair, and they often need 3D face morphable models [7, 34], or multi-view images [65, 53] to provide 3D information. Our method offers a general-purpose solution that can handle a diverse set of objects without 3D data.

3D-aware generative models. Generative image models learn to synthesize realistic 2D images [36, 81, 30, 88, 17]. However, the original formulations do not account for the 3D nature of our visual world, making 3D manipulation difficult. Recently, several methods have integrated implicit scene representation, volumetric rendering, and GANs into generative NeRFs [78, 13, 67, 29]. Given a sampled viewpoint, an image is rendered via volumetric rendering and fed to a discriminator. For example, EG3D [12] uses an efficient 3D representation called tri-planes, and StyleSDF [68] merges the style-based architecture and the SDF-based volume renderer. Multiple works [12, 68, 105, 37] have developed a two-stage model to generate high-resolution images. With GAN inversion methods [108, 74, 33, 70, 27, 111, 31, 4], we can utilize these 3D-aware generative models to align and blend images and produce multi-view consistent 3D effects.

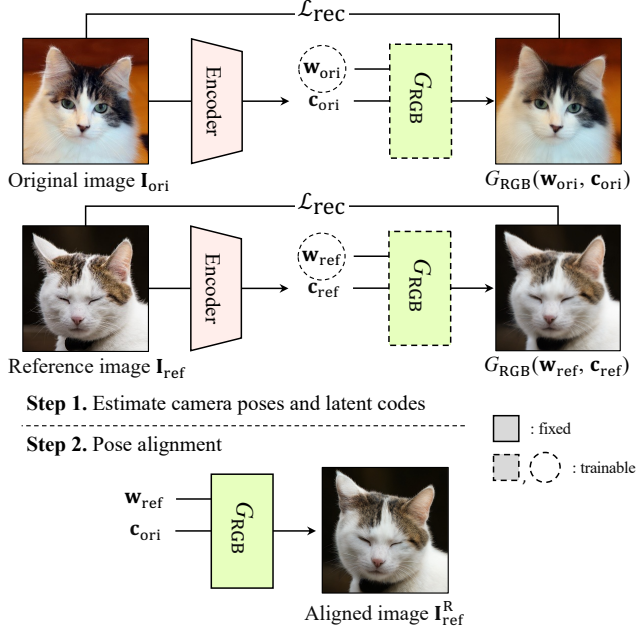


Figure 3: **3D-aware alignment**: we first use a CNN encoder to infer the camera pose of each input image. **Step 1.** Given the camera pose c , we estimate the latent code w for each input using a reconstruction loss \mathcal{L}_{rec} . **Step 2.** Given the estimated camera pose c_{ori} and latent code w_{ref} , we align the reference image to match the pose of the original image.

3D-aware image editing. Classic 3D-aware image editing methods can create 3D effects given 2D photographs [48, 19, 49]. However, they often require manual efforts to reconstruct the input’s geometry and texture. Recently, to reduce manual efforts, researchers have employed generative NeRFs for 3D-aware editing. For example, EditNeRF [62] uses separate latent codes to edit the shape and color of a NeRF object. NeRF-Editing [98] proposes to reflect geometric edits in implicit neural representations. CLIP-NeRF [89] uses a CLIP loss [73] to ensure that the edited result corresponds to the input condition. In SURF-GAN [57], they discover controllable attributes using NeRFs for training a 3D-controllable GAN. Kobayashi et al. [56] enable editing via semantic scene decomposition. While the above works tackle various image editing tasks, we focus on a different task – image blending, which requires both alignment and harmonization. Compared to previous image blending methods, our method addresses blending in a 3D-aware manner.

3. Method

We aim to perform 3D-aware image blending using only 2D images, with target masks from users for both original and reference images. Our method consists of two stages: 3D-aware alignment and 3D-aware blending. Before we blend, we first align the pair of images regarding the pose.

In Section 3.1, we describe *pose alignment* for entire objects and *local alignment* for target regions. Then, we apply the 3D-aware blending method in the generative NeRF’s latent space in Section 3.2. A variation of our blending method is illustrated in Section 3.3. We combine Poisson blending with our method to achieve near-perfect background preservation. We use EG3D [12] as our backbone, although other 3D-aware generative models, such as StyleSDF [68], can also be applied. See Appendix E for more details.

3.1. 3D-aware alignment

Pose alignment is a requisite process of our blending method, as slight pose misalignment of two images can severely degrade blending quality as shown in Figure 1. To match the reference image I_{ref} to the pose of the original image I_{ori} , we use a generative NeRF G to estimate the camera pose c and the latent code w of each image. In *Step 1* in Figure 3, we first train and freeze a CNN encoder (*i.e.*, pose estimation network) to predict the camera poses of input images. During training, we can generate a large number of pairs of camera poses and images using generative NeRF and train the encoder E using a pose reconstruction loss $\mathcal{L}_{\text{pose}}$ as follows:

$$\mathcal{L}_{\text{pose}} = \mathbb{E}_{w, c} \|c - E(G_{\text{RGB}}(w, c))\|_1, \quad (1)$$

where G_{RGB} is an image rendering function with the generative NeRF G , and $\|\cdot\|_1$ is the L1 distance. The latent code w and camera pose c are randomly drawn.

With our trained encoder, we estimate the camera poses c_{ori} and c_{ref} (defined as Euler angles $c \in SO(3)$) of the original and reference images, respectively. Given the estimated camera poses, we project input images I_{ori} and I_{ref} to the latent codes w_{ori} and w_{ref} using Pivotal Tuning Inversion (PTI) [74]. We optimize the latent code w using the reconstruction loss \mathcal{L}_{rec} as follows:

$$\mathcal{L}_{\text{rec}} = \|I - G_{\text{RGB}}(w, c)\|_1 + \mathcal{L}_{\text{LPIPS}}(I, G_{\text{RGB}}(w, c)), \quad (2)$$

where $\mathcal{L}_{\text{LPIPS}}$ is a learned perceptual image patch similarity (LPIPS) [103] loss. For more accurate inversion, we fine-tune the generator G . Inversion details are described in Appendix B. Finally, as shown in *Step 2* of Figure 3, we can align the reference image as follows:

$$I_{\text{ref}}^{\text{R}} = G_{\text{RGB}}(w_{\text{ref}}, c_{\text{ori}}). \quad (3)$$

While *pose alignment* can align two entire objects, further alignment in editing regions may still be necessary due to variations in scale and translation between object instances. To align target editing regions such as the face, eyes, and ears, we can further employ *local alignment* in the loosely aligned dataset (AFHQv2). The Iterative Closest Point (ICP) algorithm [5, 20] is applied to meshes, which can adjust their scale and translation. For further details, please refer to Appendix C.

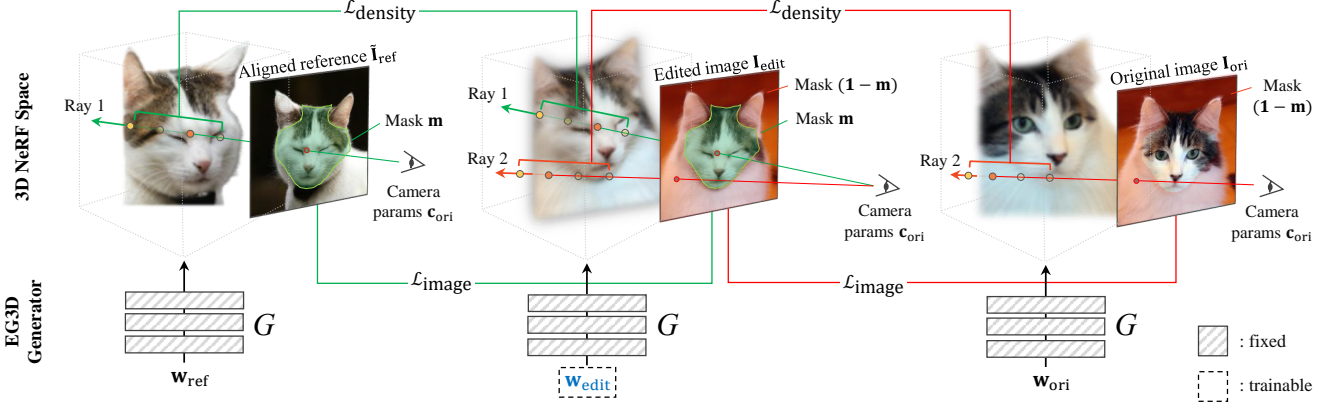


Figure 4: **Our 3D-aware blending pipeline.** We employ density-blending loss ($\mathcal{L}_{\text{density}}$) in the volume density of 3D NeRF space, as well as the image-blending loss ($\mathcal{L}_{\text{image}}$) in 2D image space. **Green rays** pass through the interior of the mask (\mathbf{m}) and **red rays** pass through the exterior of the mask ($1 - \mathbf{m}$). $\mathcal{L}_{\text{image}}$ and $\mathcal{L}_{\text{density}}$ are used to optimize the latent code \mathbf{w}_{edit} to generate the well-blended image \mathbf{I}_{edit} .

3.2. 3D-aware blending

We aim to find the best latent code \mathbf{w}_{edit} to synthesize a seamless and natural output. To achieve this goal, we exploit both 2D pixel constraints (RGB value) and 3D geometric constraints (volume density). With the proposed image-blending and density-blending losses, we optimize the latent code \mathbf{w}_{edit} , by matching the foreground with the reference and the background with the original. Algorithm details such as optimization latent spaces and input masks can be found in Appendix D.

Image-blending algorithms are often designed to match the color and details of the original image (*i.e.*, background) while preserving the structure of the reference image (*i.e.*, foreground) [72]. As shown in Figure 4, our image-blending loss matches the color and perceptual similarity of the original image using a combination of L1 and LPIPS [103], while matching the reference image’s details using LPIPS loss alone. L1 loss in the reference can lead to overfitting to the pixel space. Let \mathbf{I}_{edit} be the rendered image from the latent code \mathbf{w}_{edit} . We define the image-blending loss as follows:

$$\begin{aligned} \mathcal{L}_{\text{image}} = & \|(\mathbf{1} - \mathbf{m}) \circ \mathbf{I}_{\text{edit}} - (\mathbf{1} - \mathbf{m}) \circ \mathbf{I}_{\text{ori}}\|_1 \\ & + \lambda_1 \mathcal{L}_{\text{LPIPS}}((\mathbf{1} - \mathbf{m}) \circ \mathbf{I}_{\text{edit}}, (\mathbf{1} - \mathbf{m}) \circ \mathbf{I}_{\text{ori}}) \\ & + \lambda_2 \mathcal{L}_{\text{LPIPS}}(\mathbf{m} \circ \mathbf{I}_{\text{edit}}, \mathbf{m} \circ \mathbf{I}_{\text{ref}}), \end{aligned} \quad (4)$$

where \circ denotes element-wise multiplication. Here, λ_1 and λ_2 balance each loss term.

Density-blending is our key component in 3D-aware image blending. If we use only image-blending loss, the blending result easily falls blurry and may not reflect the reference object correctly. Especially, a highly structured object such

as hair is hard to be blended in the 3D NeRF space without volume density, as shown in Figure 8. By representing each image as a NeRF instance, we can calculate the density σ of a given 3D location $\mathbf{x} \in \mathbb{R}^3$. Let \mathcal{R}_{ref} and \mathcal{R}_{ori} be the set of rays \mathbf{r} passing through the interior and exterior of the target mask \mathbf{m} , respectively. For the 3D sample points along the rays \mathcal{R}_{ref} , we aim to match the density field between the reference and our output result, as shown as the sample points in a green ray in Figure 4. For 3D sample points in \mathcal{R}_{ori} , we also match the density field between the original and the result, as shown as the sample points in a red ray in Figure 4. Let $G_{\sigma}(\mathbf{w}; \mathbf{x})$ be the density of a given 3D point \mathbf{x} with a given latent code \mathbf{w} . Our density-blending loss can be formulated as follows:

$$\begin{aligned} \mathcal{L}_{\text{density}} = & \sum_{\mathbf{r} \in \mathcal{R}_{\text{ref}}} \sum_{\mathbf{x} \in \mathbf{r}} \|G_{\sigma}(\mathbf{w}_{\text{edit}}; \mathbf{x}) - G_{\sigma}(\mathbf{w}_{\text{ref}}; \mathbf{x})\|_1 \\ & + \sum_{\mathbf{r} \in \mathcal{R}_{\text{ori}}} \sum_{\mathbf{x} \in \mathbf{r}} \|G_{\sigma}(\mathbf{w}_{\text{edit}}; \mathbf{x}) - G_{\sigma}(\mathbf{w}_{\text{ori}}; \mathbf{x})\|_1. \end{aligned} \quad (5)$$

Our final objective function includes both image-blending loss and density-blending loss:

$$\mathcal{L} = \lambda \mathcal{L}_{\text{image}} + \mathcal{L}_{\text{density}}, \quad (6)$$

where λ is the hyperparameter that controls the contribution of the image-blending loss. If our user wants to blend the shape together without reflecting the color of reference, λ_2 in Eqn. 4 is set to 0. Otherwise, we can set λ_2 to a positive number to reflect the reference image’s color and geometry as shown in Figure 9.

3.3. Combining with Poisson blending

While our method produces high-quality blending results, incorporating Poisson blending [72] further improves the

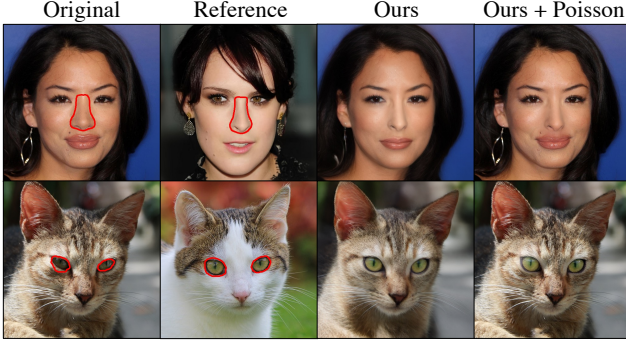


Figure 5: **Ours with Poisson blending** [72]. Ours shows satisfying blending results but a lack of preservation in details. In the first row, the earring is missing in our method. The high-frequency details such as hair and fur are less pronounced in our method alone compared to when it is combined with Poisson blending.

preservation of the original image details. Figure 5 shows the effect of Poisson blending with our method. We perform Poisson blending between the original image and the blended image generated by our 3D-aware blending method. Our blending method is modified in two ways. 1) In the initial blending stage, we only preserve the area around the border of the mask instead of all parts of the original image, as we can directly use the original image in the Poisson blending stage. We can reduce the number of iterations from 200 to 100, as improved faithfulness scores are easily achieved; see mL_2 and $LPIPS_m$ in Tables 1 and 2. 2) Instead of using the latent code of the original image w_{ori} as the initial value of w_{edit} , we use the latent code of the reference image w_{ref} . This allows us to instantly reflect the identity of the reference image and only optimize w_{edit} to reconstruct a small region near the mask boundary of the original image. Note that this is an optional choice, as our method *without* Poisson blending has already outperformed all the baselines, as shown in Tables 3 and 4.

4. Experiments

In this section, we show the advantage of our full method over several existing methods and ablated baselines. In Section 4.1, we describe our experimental settings, including baselines, datasets, and evaluation metrics. In Section 4.2, we show both quantitative and qualitative comparisons. In addition to the automatic quantitative metrics, our user study shows that our method is preferred over baselines regarding photorealism. In Section 4.3, we analyze the effectiveness of each module via ablation studies. Lastly, Section 4.4 shows useful by-products of our method, such as generating multi-view images and controlling the color and geometry disentanglement. We provide experimental details in Appendix A and runtime of our method in Appendix B, C, and D.

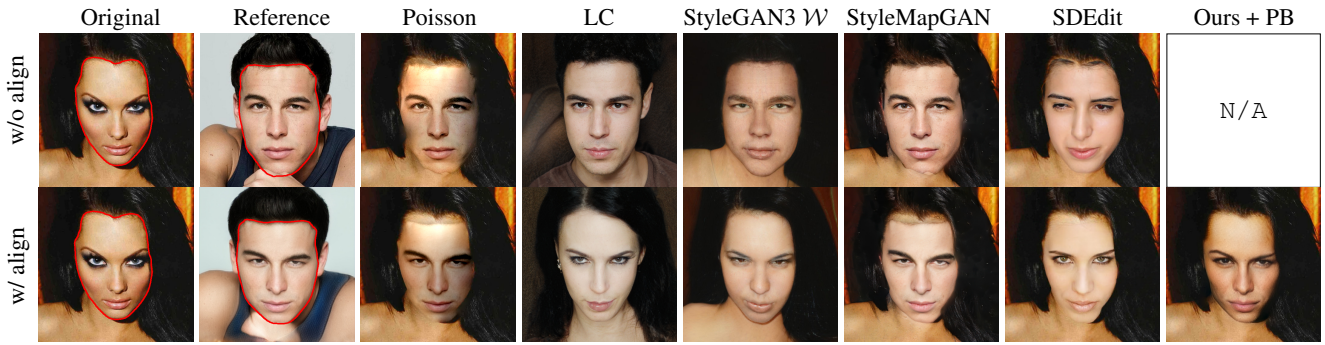
4.1. Experimental setup

Baselines. We compare our method with various image blending methods using only 2D input images. For classic methods, we run Poisson blending [72], a widely-used gradient-domain editing method. We also compare with several recent learning-based methods [11, 50, 45, 64]. Latent Composition [11] utilizes the compositionality in GANs by finding the latent code of the roughly collaged inputs on the manifold of the generator. StyleMapGAN [50] proposes the spatial latent space for GANs to enable local parts blending by mixing the spatial latent codes. Recently, Karras *et al.* [45] proposed StyleGAN3, which provides rotation equivariance. Therefore, we additionally show their blending results by finding the latent code of the composited inputs on the StyleGAN3-R manifold. Both \mathcal{W} and $\mathcal{W}+$ of StyleGAN3-R latent spaces are tested. SDEdit [64] is a diffusion-based blending method that produces a natural-looking result by denoising the corrupted image of a composite image.

Datasets. We use FFHQ [46] and AFHQv2-Cat datasets [21] for model training. We use pose alignment for both datasets and apply further local alignment to the loosely aligned dataset (AFHQ).

To test blending performance, we use CelebA-HQ [44] for the FFHQ-trained models and AFHQv2-Cat test sets for the AFHQ-trained models. We randomly select 250 pairs of images from each dataset for an original and reference image. We also create a target mask for each pair to automatically simulate a user input using pretrained semantic segmentation networks [96, 104, 16]. We blend 5 and 3 semantic parts in each pair of images for CelebA-HQ and AFHQ, respectively. The total number of blended images in each method is 1,250 (CelebA-HQ) and 750 (AFHQv2-Cat). We also include results on ShapeNet-Car dataset [14] to show that our method works well for non-facial data.

Evaluation metrics. For evaluation metrics, we use masked L_2 , masked LPIPS [103] and Kernel Inception Score (KID) [6]. Masked L_2 (mL_2) is the L_2 distance between the original image and the blended image on the exterior of the mask, measuring the preservation of non-target areas of the original image. Unlike background regions, a pixel-wise loss is too strict for the target area changed during blending. We measure the perceptual similarity metric (LPIPS) [103] for the blended regions, which are called masked LPIPS ($LPIPS_m$) used in previous methods [41, 64]. Kernel Inception Score (KID) [6] is widely used to quantify the realism of the generated images regarding the real data distribution. We compute KID between blended images and the training dataset using the `clean-fid` library [71].



Method	w/o align (baseline only)			w/ 3D-aware align		
	KID ↓	LPIPS _m ↓	mL ₂ ↓	KID ↓	LPIPS _m ↓	mL ₂ ↓
Poisson Blending [72]	<u>0.006</u>	0.4203	0.0069	<u>0.005</u>	0.2355	0.0051
Latent Composition [11]	0.012	0.4735	0.0388	0.012	0.4487	0.0321
StyleGAN3 \mathcal{W} [45]	0.016	0.4379	0.0353	0.017	0.3921	0.0307
StyleGAN3 $\mathcal{W}+$ [45]	0.025	0.4634	0.0462	0.023	0.4086	0.0391
StyleMapGAN (32 × 32) [50]	0.007	0.3792	0.0118	0.006	<u>0.1989</u>	0.0045
SDEdit [64]	0.011	0.3857	0.0076	0.008	0.3427	0.0003
Ours	0.013	<u>0.2046</u>	<u>0.0050</u>	0.013	0.2046	0.0050
Ours + Poisson Blending	0.002	0.1883	0.0007	0.002	0.1883	<u>0.0007</u>

Table 1: **Comparison with baselines in the CelebA-HQ test set.** The first and second rows of the *figure* show the blending results without and with our 3D-aware alignment, respectively. Metric scores on the left side of the *table* show the results without alignment. We apply our 3D-aware alignment to the baselines on the right side of the table. Lower scores denote better performance in all metrics. The best and second-best scores are bold and underlined. Our method outperforms baselines in all metrics. LC and PB stand for Latent Composition [11] and Poisson Blending [72], respectively. Note that our method always operates 3D-aware alignment, as it is an integral part of our algorithm.

User study. To further examine the effectiveness of our 3D-aware blending method, we conduct a user study for photorealism. Our target goal is to edit the original image, so we exclude baselines that show highly flawed preservation of the original image. Human evaluates pairwise comparison of blended images between our method and one of the baselines. The user selects more real-looking images. We collect 5,000 comparison results via Amazon Mechanical Turk (MTurk). See Appendix F for more details.

4.2. Comparison with baselines

Here we compare our method with baselines in two variations. In the *w/o align* setting, we do not apply our 3D-aware alignment to baselines. In the *w/ align* setting, we align the reference image with our 3D-aware alignment. This experiment demonstrates the effectiveness of our proposed method. 1) Our alignment method consistently improves all baselines in all evaluation metrics: KID, LPIPS_m, and masked L_2 . 2) Our 3D-aware blending method outperforms all baselines, including those that use our alignment method. We also report the combination of our method and Poisson blending to achieve better background preservation, as the perfect inversion is still hard to be achieved in GAN-based methods.

Table 1 shows comparison results in CelebA-HQ. The left side of the table includes all the baselines without our 3D-aware alignment. All metrics are worse than the right side of the table (*w/ align*). This result reveals that alignment between the original and reference image affects overall editing performance. Table 2 shows comparison results in AFHQv2-Cat. It shows the same tendency as Table 1. More comparison results are included in Appendix I.

Our method performs well regarding all metrics. Combined with Poisson blending, our method outperforms all baselines. Poisson blending and StyleMapGAN (16 × 16, 32 × 32) show great faithfulness to the input images but suffer from artifacts. Latent Composition, StyleMapGAN (8 × 8), and StyleGAN3 \mathcal{W} produce realistic results but far from the input images. The identities of the original and reference images have changed, which is reflected by a worse LPIPS_m and mL₂. SDEdit fails to reflect the reference image and shows worse LPIPS_m. StyleGAN3 $\mathcal{W}+$ often shows entirely collapsed images. Our method preserves the identity of the original image and reflects the reference image well while producing realistic outputs.



Method	w/o align (baseline only)			w/ 3D-aware align		
	KID ↓	LPIPS _m ↓	mL ₂ ↓	KID ↓	LPIPS _m ↓	mL ₂ ↓
Poisson Blending [72]	0.002	0.4956	<u>0.0024</u>	0.002	<u>0.2656</u>	0.0004
StyleGAN3 \mathcal{W} [45]	0.006	0.4588	0.0316	0.006	0.3802	0.0268
StyleGAN3 $\mathcal{W}+$ [45]	0.014	0.4941	0.0298	0.013	0.3903	0.0236
StyleMapGAN (8 × 8) [50]	0.013	0.4840	0.0574	0.013	0.3221	0.0526
StyleMapGAN (16 × 16) [50]	0.006	0.4746	0.0225	0.004	0.2707	0.0160
Ours	0.005	<u>0.2739</u>	0.0073	0.005	0.2739	0.0073
Ours + Poisson Blending	0.002	0.2229	0.0013	0.002	0.2229	<u>0.0013</u>

Table 2: Comparison with baselines in the AFHQv2-Cat test set. Formats of the figure and table are the same as Table 1.

Method	Ours		Ours + Poisson Blending	
	w/o	w/ align	w/o	w/ align
Poisson [72]	79.9%	59.9%	80.9% (+1.0)	67.5% (+7.6)
StyleMap [50]	72.3%	62.0%	75.4% (+3.1)	66.3% (+4.3)
SDEdit [64]	61.0%	55.7%	61.1% (+0.1)	50.2% (-5.5)

Table 3: User study in CelebA-HQ regarding the photorealism of the blended image. The percentage denotes how often MTurk workers prefer our method to each baseline in pairwise comparison. Values larger than 50% mean ours outperforms the baseline. Our method, both with and without Poisson blending, outperforms all baselines even if we improve the baselines using our 3D-aware alignment. Incorporating Poisson blending further enhances the realism score of our method as shown in green numbers.

Method	Ours		Ours + Poisson Blending	
	w/o	w/ align	w/o	w/ align
Poisson [72]	91.2%	76.2%	87.6% (-3.6)	82.8% (+6.6)
StyleMap [50]	91.0%	82.3%	91.6% (+0.6)	83.6% (+1.3)

Table 4: User study in AFHQv2-Cat regarding the photorealism of the blended image. All details are the same as in Table 3. Our approach surpasses all baselines, and the incorporation of Poisson blending further improves the realism score.

User study. We note that KID has a high correlation with background preservation. Unfortunately, it fails to capture the boundary artifacts and foreground image quality, espe-

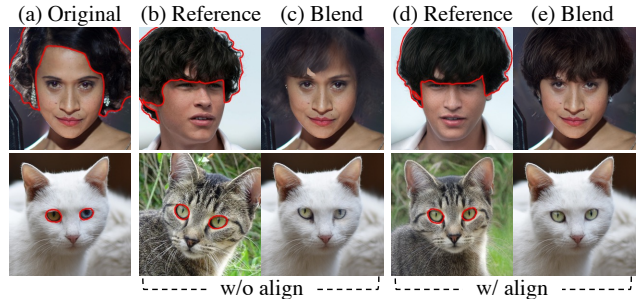


Figure 6: The effect of our 3D-aware alignment. Aligned reference images (d) have the same pose as the original images (a). With our alignment, blending results (e) look more realistic and reflect the reference well than those without alignment (c).

cially for small foregrounds. To further evaluate the realism score of results, we conduct a human perception study for our method and other baselines, which shows great preservation scores mL₂. As shown in Tables 3 and 4, MTurk participants prefer our method to other baselines regarding the photorealism of the results. Our method, as well as the combination of ours with Poisson blending, outperforms the baselines. SDEdit with our 3D-aware alignment shows a comparable realism score with ours, but it can not reflect the reference well, as reflected in worse LPIPS_m score in Table 1. Similar to Tables 1 and 2, Mturk participants prefer baselines with alignment to their unaligned counterparts.



Figure 7: Ablation study of Poisson blending (PB) in baselines. Despite combining Poisson blending with the baselines, StyleMapGAN still generates artifacts, and other baselines fail to preserve the identity of the reference. Our method with Poisson blending keeps the original image intact while accurately reflecting the reference image.

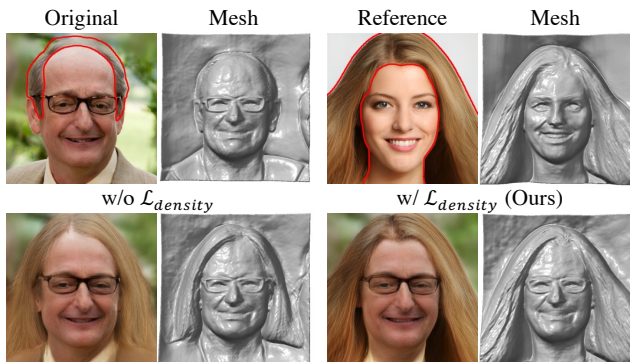


Figure 8: **The effect of our density-blending loss.** Without the loss, 3D information is not considered, resulting in inaccurate blending in 3D space. In the bottom left figure, the hair mesh is not properly reflected without the density-blending loss, resulting in inaccurate blending and missing fine details.

4.3. Ablation study

3D-aware alignment is an essential part of image blending. As discussed in Section 4.2, our alignment provides consistent improvements in all baseline methods. Moreover, it plays a crucial role in our blending approach. Figure 6 shows the importance of 3D-aware alignment, where the lack of alignment in the reference images result in degraded blending results (Figure 6c). Specifically, the woman’s hair appears blurry, and the size of the cat’s eyes looks different. Aligned reference images can generate realistic blending results (Figure 6e) in our 3D-aware blending method.

Density-blending loss gives rich 3D signals in the blending procedure. Section 3.2 explains how we can exploit volume density fields in blending. Delicate geometric structures, such as hair, can not be easily blended without awareness of 3D information. Figure 8 shows an ablation study of our density-blending loss. In the bottom left, the hair looks



Figure 9: **Color-geometry disentanglement with our model.** We can adjust the reflection of the reference image’s color by adjusting the weight λ_2 in the image-blending loss. Without image blending loss on reference, we can focus on object shapes, as shown in the rightmost column.

blurry in the blended image, and the mesh of the result shows shorter hair than that in the reference image. In the bottom right, the well-blended image and corresponding mesh show that our density-blending loss contributes to capturing the highly structured object in blending. We provide an additional ablation study using StyleSDF [68] in Appendix E.

Combination with Poisson blending. In Tables 1 and 2, we report the combination of our method and Poisson blending. It shows Poisson blending further enhances the performance of our method in all automatic metrics: KID, mL_2 , and $LPIPS_m$. In the realism score of human perception, ours with Poisson blending enhance the score, as shown in green numbers of Tables 3 and 4. However, combining Poisson blending with each baseline does not have meaningful benefits, as shown in Figure 7. Baselines still show artifacts or fail to reflect the identity of the reference.

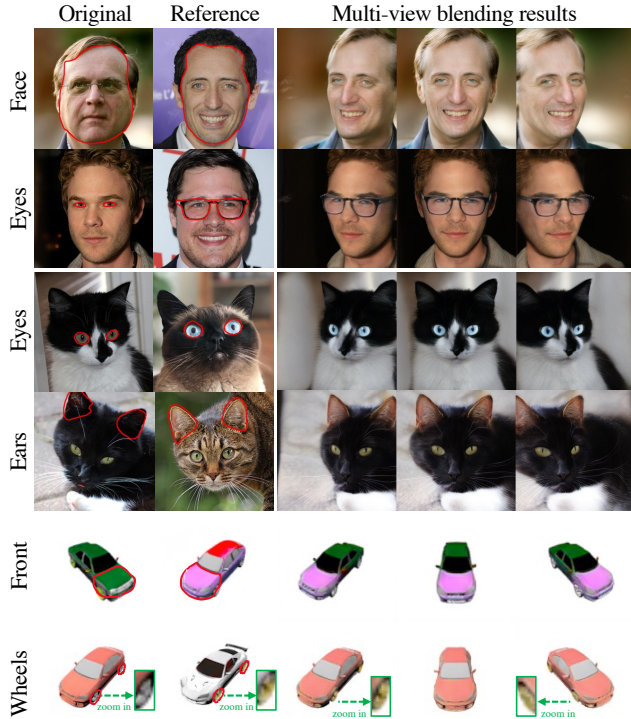


Figure 10: **Multi-view blending results** in various datasets: CelebA-HQ, AFHQv2-Cat, and ShapeNet-Car. Since we optimize the latent code of the generative NeRF, we can synthesis images of the blended object in different poses through the generative NeRF.

4.4. Additional advantages of NeRF-based blending

In addition to increasing blending quality, our 3D-aware method enables additional capacity: color-geometry disentanglement and multi-view consistent blending. As shown in Figure 9, we can control the influence of color in blending. The results with $\mathcal{L}_{\text{image}}$ have a redder color than the results without the loss. If we remove or assign a lower weight to the image-blending loss on reference (λ_2 in Eqn. 4), we can reflect the geometry of the reference object more than the color. In contrast, we can reflect colors better if we give a larger weight to λ_2 . Note that we always use the image-blending loss on the original image to preserve it better. An additional ablation study using StyleSDF is included in Appendix E.

A key component of generative NeRFs is multi-view consistent generation. After applying the blending procedure described in Section 3.2, we have an optimized latent code \mathbf{w}_{edit} . Generative NeRF can synthesize a novel view blended image using \mathbf{w}_{edit} and a target camera pose. Figure 10 shows the multi-view consistent blending results in CelebA-HQ, AFHQv2-Cat, and ShapeNet-Car [14]. We provide more multi-view results for EG3D and StyleSDF in Appendix I.

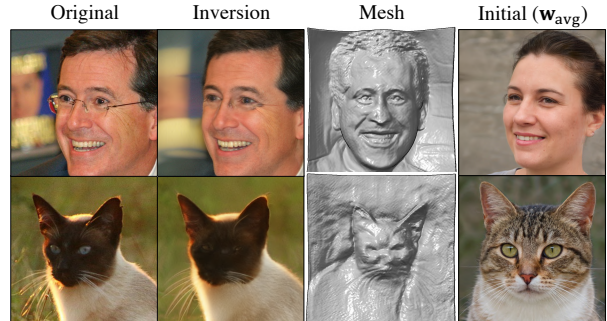


Figure 11: Failure cases of inversion. If an input image has a large variance in scale relative to the mean face or the estimated pose from the encoder is not valid, inversion sometimes fails. The first row shows a failure to reconstruct eyeglasses, and the second row shows a crushed face of a cat in the reconstructed image and mesh.

5. Discussion and Limitations

Our method exploits the capability of NeRFs to align and blend images in a 3D-aware manner only with a collection of 2D images. Our 3D-aware alignment boosts the quality of existing 2D baselines. 3D-aware blending exceeds improved 2D baselines with our alignment method and shows additional advantages such as color-geometry disentanglement and multi-view consistent blending. We hope our approach paves the road to 3D-aware blending. Recently, 3DGP [80] presents a 3D-aware GAN, handling non-alignable scenes captured from arbitrary camera poses in real-world environments. Since our approach relies solely on a pre-trained generator, it can be readily extended to blend unaligned multi-category datasets such as ImageNet [28].

Despite improvements over existing blending baselines, our method depends on GAN inversion, which is a bottleneck of the overall performance regarding quality and speed. Figure 11 shows the inversion process can sometimes fail to accurately reconstruct the input image. We cannot obtain an acceptable inversion result if an input image is far from the average face generated from the mean latent code \mathbf{w}_{avg} . We also note the camera pose inferred by our encoder should not be overly inaccurate. Currently, the problem is being addressed by combining our method with Poisson blending. However, more effective solutions may be available with recent advances in 3D GAN inversion techniques [92, 55]. In the future, to enable real-time editing, we could explore training an encoder [31, 50] to blend images using our proposed loss functions.

Acknowledgments. We would like to thank Seohui Jeong, Che-Sang Park, Eric R. Chan, Junho Kim, Jung-Woo Ha, Youngjung Uh, and other NAVER AI Lab researchers for their helpful comments and sharing of materials. All experiments were conducted on NAVER Smart Machine Learning (NSML) platform [51, 82].

References

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *ICCV*, 2019. 14, 16
- [2] Aseem Agarwala, Mira Dontcheva, Maneesh Agrawala, Steven Drucker, Alex Colburn, Brian Curless, David Salesin, and Michael Cohen. Interactive digital photomontage. In *ACM SIGGRAPH*. 2004. 2
- [3] Dmitry Baranchuk, Ivan Rubachev, Andrey Voynov, Valentin Khrukov, and Artem Babenko. Label-efficient semantic segmentation with diffusion models. In *International Conference on Learning Representations (ICLR)*, 2022. 14
- [4] David Bau, Hendrik Strobelt, William Peebles, Jonas Wulff, Bolei Zhou, Jun-Yan Zhu, and Antonio Torralba. Semantic photo manipulation with a generative image prior. *arXiv preprint arXiv:2005.07727*, 2020. 2
- [5] Paul J Besl and Neil D McKay. Method for registration of 3-d shapes. In *Sensor fusion IV: control paradigms and data structures*, 1992. 3, 15, 16, 20
- [6] Mikołaj Bińkowski, Danica J. Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. In *International Conference on Learning Representations (ICLR)*, 2018. 5
- [7] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, 1999. 2
- [8] G. Bradski. The OpenCV Library. *Dr. Dobb's Journal of Software Tools*, 2000. 14
- [9] Matthew Brown, David G Lowe, et al. Recognising panoramas. In *IEEE International Conference on Computer Vision (ICCV)*, 2003. 2
- [10] Peter J Burt and Edward H Adelson. The laplacian pyramid as a compact image code. In *Readings in computer vision*. 1987. 1, 2
- [11] Lucy Chai, Jonas Wulff, and Phillip Isola. Using latent space regression to analyze and leverage compositionality in gans. In *International Conference on Learning Representations (ICLR)*, 2021. 2, 5, 6, 14, 18, 20, 22
- [12] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1, 2, 3, 14, 15, 18
- [13] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 2
- [14] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 5, 9, 14, 24
- [15] Bor-Chun Chen and Andrew Kae. Toward realistic image compositing with adversarial learning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1
- [16] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017. 5
- [17] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pre-training from pixels. In *International Conference on Machine Learning (ICML)*, 2020. 2
- [18] Tao Chen, Ming-Ming Cheng, Ping Tan, Ariel Shamir, and Shi-Min Hu. Sketch2photo: Internet image montage. *ACM Transactions on graphics (TOG)*, 2009. 2
- [19] Tao Chen, Zhe Zhu, Ariel Shamir, Shi-Min Hu, and Daniel Cohen-Or. 3-sweep: Extracting editable objects from a single photo. *ACM Transactions on graphics (TOG)*, 2013. 3
- [20] Yang Chen and Gérard Medioni. Object modelling by registration of multiple range images. *Image and vision computing*, 1992. 3, 15, 16, 20
- [21] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 5, 14, 16, 19, 20
- [22] Christopher Choy, Wei Dong, and Vladlen Koltun. Deep global registration. In *CVPR*, 2020. 20
- [23] Chaeyeon Chung, Taewoo Kim, Hyelin Nam, Seunghwan Choi, Gyojung Gu, Sunghyun Park, and Jaegul Choo. Hairfit: pose-invariant hairstyle transfer via flow-based hair alignment and semantic-region-aware inpainting. In *The British Machine Vision Conference (BMVC)*, 2021. 2
- [24] Edo Collins, Raja Bala, Bob Price, and Sabine Susstrunk. Editing in style: Uncovering the local semantics of gans. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [25] Edo Collins, Raja Bala, Bob Price, and Sabine Susstrunk. Editing in style: Uncovering the local semantics of gans. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [26] Kevin Dale, Kalyan Sunkavalli, Micah K Johnson, Daniel Vlasic, Wojciech Matusik, and Hanspeter Pfister. Video face replacement. In *ACM SIGGRAPH Asia*, 2011. 2
- [27] Giannis Daras, Wen-Sheng Chu, Abhishek Kumar, Dmitry Lagun, and Alexandros G Dimakis. Solving inverse problems with nerfgans. *arXiv preprint arXiv:2112.09061*, 2021. 2
- [28] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 9
- [29] Yu Deng, Jiaolong Yang, Jianfeng Xiang, and Xin Tong. Gram: Generative radiance manifolds for 3d-aware image generation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [30] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2021. 2, 14

- [31] Tan M Dinh, Anh Tuan Tran, Rang Nguyen, and Binh-Son Hua. Hyperinverter: Improving stylegan inversion via hypernetwork. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 9
- [32] Zeev Farbman, Gil Hoffer, Yaron Lipman, Daniel Cohen-Or, and Dani Lischinski. Coordinates for instant image cloning. *ACM Transactions on graphics (TOG)*, 2009. 2
- [33] Qianli Feng, Viraj Shah, Raghudeep Gadde, Pietro Perona, and Aleix Martinez. Near perfect gan inversion. *arXiv preprint arXiv:2202.11833*, 2022. 2
- [34] Claudio Ferrari, Stefano Berretti, Pietro Pala, and Alberto Del Bimbo. A sparse and locally coherent morphable face model for dense semantic correspondence across heterogeneous 3d faces. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2021. 2
- [35] Anna Frühstück, Ibraheem Alhashim, and Peter Wonka. Tigan: synthesis of large-scale non-homogeneous textures. *ACM Transactions on graphics (TOG)*, 2019. 2
- [36] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2014. 2
- [37] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. Stylenerf: A style-based 3d-aware generator for high-resolution image synthesis. In *International Conference on Learning Representations (ICLR)*, 2022. 1, 2
- [38] Jeffrey T Hancock and Jeremy N Bailenson. The social impact of deepfakes. *Cyberpsychology, behavior, and social networking*, 2021. 20
- [39] James Hays and Alexei A Efros. Scene completion using millions of photographs. *ACM Transactions on graphics (TOG)*, 2007. 1, 2
- [40] Berthold K.P. Horn and Brian G. Schunck. Determining optical flow. *Artificial Intelligence*, 1981. 16
- [41] Minyoung Huh, Richard Zhang, Jun-Yan Zhu, Sylvain Paris, and Aaron Hertzmann. Transforming and projecting images into class-conditional generative networks. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 17–34. Springer, 2020. 5
- [42] Jiaya Jia, Jian Sun, Chi-Keung Tang, and Heung-Yeung Shum. Drag-and-drop pasting. In *ACM Transactions on graphics (TOG)*, 2006. 2
- [43] Omer Kafri, Or Patashnik, Yuval Alaluf, and Daniel Cohen-Or. Stylefusion: A generative model for disentangling spatial segments. *arXiv preprint arXiv:2107.07437*, 2021. 20
- [44] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. In *International Conference on Learning Representations (ICLR)*, 2018. 5, 19, 20
- [45] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2021. 5, 6, 7, 14, 18
- [46] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 5, 14, 16
- [47] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 14
- [48] Kevin Karsch, Varsha Hedau, David Forsyth, and Derek Hoiem. Rendering synthetic objects into legacy photographs. *ACM Transactions on graphics (TOG)*, 2011. 3
- [49] Natasha Kholgade, Tomas Simon, Alexei Efros, and Yaser Sheikh. 3d object manipulation in a single photograph using stock 3d models. *ACM Transactions on graphics (TOG)*, 2014. 3
- [50] Hyunsu Kim, Yunjeon Choi, Junho Kim, Sungjoo Yoo, and Youngjung Uh. Exploiting spatial dimensions of latent in gan for real-time image editing. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 2, 5, 6, 7, 9, 14, 18, 20
- [51] Hanjoo Kim, Minkyu Kim, Dongjoo Seo, Jinwoong Kim, Heungseok Park, Soeun Park, Hyunwoo Jo, KyungHyun Kim, Youngil Yang, Youngkwan Kim, et al. Nsml: Meet the mlaas platform with a real-world case study. *arXiv preprint arXiv:1810.09957*, 2018. 9
- [52] Taewoo Kim, Chaeyeon Chung, Yoonseo Kim, Sunghyun Park, Kangyeol Kim, and Jaegul Choo. Style your hair: Latent optimization for pose-invariant hairstyle transfer via local-style-aware hair alignment. In *European Conference on Computer Vision (ECCV)*, 2022. 2
- [53] Taewoo Kim, Chaeyeon Chung, Sunghyun Park, Gyojung Gu, Keonmin Nam, Wonzo Choe, Jaesung Lee, and Jaegul Choo. K-hairstyle: A large-scale korean hairstyle dataset for virtual hair editing and hairstyle classification. In *IEEE International Conference on Image Processing (ICIP)*, 2021. 2
- [54] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015. 14
- [55] Jaehoon Ko, Kyusun Cho, Daewon Choi, Kwangrok Ryoo, and Seungryong Kim. 3d gan inversion with pose optimization. *WACV*, 2023. 9
- [56] Sosuke Kobayashi, Eiichi Matsumoto, and Vincent Sitzmann. Decomposing nerf for editing via feature field distillation. *arXiv preprint arXiv:2205.15585*, 2022. 3
- [57] Jeong-gi Kwak, Yuanming Li, Dongsik Yoon, Donghyeon Kim, David Han, and Hanseok Ko. Injecting 3d perception of controllable nerf-gan into stylegan for editable portrait image synthesis. In *European Conference on Computer Vision (ECCV)*, 2022. 3
- [58] Vivek Kwatra, Arno Schödl, Irfan Essa, Greg Turk, and Aaron Bobick. Graphcut textures: Image and video synthesis using graph cuts. *ACM Transactions on graphics (TOG)*, 2003. 1, 2
- [59] Jean-François Lalonde, Derek Hoiem, Alexei A Efros, Carsten Rother, John Winn, and Antonio Criminisi. Photo clip art. *ACM Transactions on graphics (TOG)*, 2007. 2

- [60] Anat Levin, Dani Lischinski, and Yair Weiss. A closed-form solution to natural image matting. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2007. [2](#)
- [61] Chen-Hsuan Lin, Ersin Yumer, Oliver Wang, Eli Shechtman, and Simon Lucey. St-gan: Spatial transformer generative adversarial networks for image compositing. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. [1](#)
- [62] Steven Liu, Xiuming Zhang, Zhoutong Zhang, Richard Zhang, Jun-Yan Zhu, and Bryan Russell. Editing conditional radiance fields. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. [3](#)
- [63] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. *ACM Transactions on graphics (TOG)*, 1987. [15](#)
- [64] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. In *International Conference on Learning Representations (ICLR)*, 2021. [1](#), [2](#), [5](#), [6](#), [7](#), [14](#), [18](#), [19](#)
- [65] Arsha Nagrani, Joon Son Chung, and Andrew Senior. Voxceleb: a large-scale speaker identification dataset. In *INTERSPEECH*, 2017. [2](#)
- [66] Thanh Thi Nguyen, Quoc Viet Hung Nguyen, Dung Tien Nguyen, Duc Thanh Nguyen, Thien Huynh-The, Saeid Nahavandi, Thanh Tam Nguyen, Quoc-Viet Pham, and Cuong M Nguyen. Deep learning for deepfakes creation and detection: A survey. *Computer Vision and Image Understanding*, 2022. [2](#)
- [67] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. [2](#)
- [68] Roy Or-El, Xuan Luo, Mengyi Shan, Eli Shechtman, Jeong Joon Park, and Ira Kemelmacher-Shlizerman. Stylesdf: High-resolution 3d-consistent image and geometry generation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. [2](#), [3](#), [8](#), [15](#), [18](#), [24](#)
- [69] Xingang Pan, Xudong Xu, Chen Change Loy, Christian Theobalt, and Bo Dai. A shading-guided generative implicit model for shape-accurate 3d-aware image synthesis. *Conference on Neural Information Processing Systems (NeurIPS)*, 2021. [1](#)
- [70] Gaurav Parmar, Yijun Li, Jingwan Lu, Richard Zhang, Jun-Yan Zhu, and Krishna Kumar Singh. Spatially-adaptive multilayer selection for gan inversion and editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11399–11409, 2022. [2](#)
- [71] Gaurav Parmar, Richard Zhang, and Jun-Yan Zhu. On aliased resizing and surprising subtleties in gan evaluation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11410–11420, 2022. [5](#)
- [72] Patrick Pérez, Michel Gangnet, and Andrew Blake. Poisson image editing. In *ACM SIGGRAPH*, 2003. [1](#), [2](#), [4](#), [5](#), [6](#), [7](#), [14](#), [18](#), [22](#)
- [73] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning (ICML)*, 2021. [3](#)
- [74] Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel Cohen-Or. Pivotal tuning for latent-based editing of real images. *ACM Transactions on graphics (TOG)*, 2022. [2](#), [3](#), [15](#)
- [75] Carsten Rother, Vladimir Kolmogorov, and Andrew Blake. ” grabcut” interactive foreground extraction using iterated graph cuts. *ACM Transactions on graphics (TOG)*, 2004. [2](#)
- [76] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2016. [15](#)
- [77] Othman Sbai, Camille Couprie, and Mathieu Aubry. Surprising image compositions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshop*, 2021. [1](#)
- [78] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *Conference on Neural Information Processing Systems (NeurIPS)*, 2020. [1](#), [2](#)
- [79] Yichun Shi, Xiao Yang, Yangyue Wan, and Xiaohui Shen. Semanticstylegan: Learning compositional generative priors for controllable image synthesis and editing. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. [2](#)
- [80] Ivan Skorokhodov, Aliaksandr Siarohin, Yinghao Xu, Jian Ren, Hsin-Ying Lee, Peter Wonka, and Sergey Tulyakov. 3d generation on imagenet. In *ICLR*, 2023. [9](#)
- [81] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations (ICLR)*, 2021. [2](#)
- [82] Nako Sung, Minkyu Kim, Hyunwoo Jo, Youngil Yang, Jingwoong Kim, Leonard Lausen, Youngkwon Kim, Gayoung Lee, Donghyun Kwak, Jung-Woo Ha, et al. Nsml: A machine learning platform that enables you to focus on your models. *arXiv preprint arXiv:1712.05902*, 2017. [9](#)
- [83] Kalyan Sunkavalli, Micah K Johnson, Wojciech Matusik, and Hanspeter Pfister. Multi-scale image harmonization. *ACM Transactions on Graphics (TOG)*, 29(4):1–10, 2010. [2](#)
- [84] Ryohei Suzuki, Masanori Koyama, Takeru Miyato, Taizan Yonetsuji, and Huachun Zhu. Spatially controllable image synthesis with internal representation collaging. In *arXiv preprint arXiv:1811.10153*, 2018. [2](#)
- [85] Richard Szeliski, Matthew Uyttendaele, and Drew Steedly. Fast poisson blending using multi-splines. In *IEEE International Conference on Computational Photography (ICCP)*, 2011. [2](#)
- [86] Michael W Tao, Micah K Johnson, and Sylvain Paris. Error-tolerant image compositing. In *European Conference on Computer Vision (ECCV)*, 2010. [2](#)

- [87] Matthew Uyttendaele, Ashley Eden, and Richard Skeliski. Eliminating ghosting and exposure artifacts in image mosaics. In *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2001. 1
- [88] Aaron Van den Oord, Nal Kalchbrenner, Lasse Espeholt, Oriol Vinyals, Alex Graves, et al. Conditional image generation with pixelcnn decoders. *Conference on Neural Information Processing Systems (NeurIPS)*, 2016. 2
- [89] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Clip-nerf: Text-and-image driven manipulation of neural radiance fields. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [90] Yue Wang and Justin M. Solomon. Deep closest point: Learning representations for point cloud registration. In *IEEE International Conference on Computer Vision (ICCV)*, October 2019. 20
- [91] Huikai Wu, Shuai Zheng, Junge Zhang, and Kaiqi Huang. Gp-gan: Towards realistic high-resolution image blending. In *ACM international conference on multimedia (ACM-MM)*, 2019. 1, 2
- [92] Jiaxin Xie, Hao Ouyang, Jintan Piao, Chenyang Lei, and Qifeng Chen. High-fidelity 3d gan inversion by pseudo-multi-view optimization. In *CVPR*, 2023. 9
- [93] Yangyang Xu, Bailin Deng, Junle Wang, Yanqing Jing, Jia Pan, and Shengfeng He. High-resolution face swapping via latent semantics disentanglement. In *CVPR*, 2022. 2
- [94] Su Xue, Aseem Agarwala, Julie Dorsey, and Holly Rushmeier. Understanding and improving the realism of image composites. *ACM Transactions on graphics (TOG)*, 31(4):1–10, 2012. 2
- [95] Fei Yang, Jue Wang, Eli Shechtman, Lubomir Bourdev, and Dimitri Metaxas. Expression flow for 3d-aware face component transfer. In *ACM SIGGRAPH*, 2011. 2
- [96] Changqian Yu, Jingbo Wang, Chao Peng, Changxin Gao, Gang Yu, and Nong Sang. Bisenet: Bilateral segmentation network for real-time semantic segmentation. In *European Conference on Computer Vision (ECCV)*, 2018. 5
- [97] Ning Yu, Larry S Davis, and Mario Fritz. Attributing fake images to gans: Learning and analyzing gan fingerprints. In *IEEE International Conference on Computer Vision (ICCV)*, 2019. 20
- [98] Yu-Jie Yuan, Yang-Tian Sun, Yu-Kun Lai, Yuewen Ma, Rongfei Jia, and Lin Gao. Nerf-editing: geometry editing of neural radiance fields. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [99] Fangneng Zhan, Hongyuan Zhu, and Shijian Lu. Spatial fusion gan for image synthesis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1
- [100] He Zhang, Jianming Zhang, Federico Perazzi, Zhe Lin, and Vishal M Patel. Deep image compositing. In *Winter Conference on Applications of Computer Vision (WACV)*, 2021. 2
- [101] Lingzhi Zhang, Tarmily Wen, and Jianbo Shi. Deep image blending. In *Winter Conference on Applications of Computer Vision (WACV)*, 2020. 2
- [102] Richard Zhang, Phillip Isola, and Alexei A Efros. Colorful image colorization. In *European Conference on Computer Vision (ECCV)*, 2016. 19
- [103] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 3, 4, 5, 14, 15
- [104] Yuxuan Zhang, Huan Ling, Jun Gao, Kangxue Yin, Jean-Francois Lafleche, Adela Barriuso, Antonio Torralba, and Sanja Fidler. Datasetgan: Efficient labeled data factory with minimal human effort. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 5
- [105] Peng Zhou, Lingxi Xie, Bingbing Ni, and Qi Tian. Cips-3d: A 3d-aware generator of gans based on conditionally-independent pixel synthesis. *arXiv preprint arXiv:2110.09788*, 2021. 1, 2
- [106] Jiapeng Zhu, Yujun Shen, Deli Zhao, and Bolei Zhou. In-domain gan inversion for real image editing. In *European Conference on Computer Vision (ECCV)*, 2020. 2
- [107] Jun-Yan Zhu, Philipp Krahenbuhl, Eli Shechtman, and Alexei A Efros. Learning a discriminative model for the perception of realism in composite images. In *IEEE International Conference on Computer Vision (ICCV)*, 2015. 1
- [108] Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman, and Alexei A. Efros. Generative visual manipulation on the natural image manifold. In *European Conference on Computer Vision (ECCV)*, 2016. 2, 16
- [109] Peihao Zhu, Rameen Abdal, John Femiani, and Peter Wonka. Barbershop: Gan-based image compositing using segmentation masks. In *ACM Transactions on graphics (TOG)*, 2021. 1, 2, 20
- [110] Peihao Zhu, Rameen Abdal, John Femiani, and Peter Wonka. Hairnet: Hairstyle transfer with pose changes. In *European Conference on Computer Vision (ECCV)*, 2022. 1, 2
- [111] Peihao Zhu, Rameen Abdal, Yipeng Qin, John Femiani, and Peter Wonka. Improved stylegan embedding: Where are the good latents? *arXiv preprint arXiv:2012.09036*, 2020. 2

Overview of Appendix.

- Please refer to the code, data, and results on our website: <https://blandocs.github.io/blenderf>.
- Experimental details are described in Appendix A.
- Details of inversion are described in Appendix B.
- Details of local alignment are described in Appendix C with an extra user study.
- Details of 3D-aware blending are described in Appendix D.
- 3D-aware blending in StyleSDF is described in Appendix E. Our method can be applied to Signed Distance Fields (SDF) beyond NeRFs.
- Details of user studies are described in Appendix F.
- Failure cases are described in Appendix G.
- Societal impact is discussed in Appendix H.
- Additional qualitative results are in Appendix I.

A. Experimental details

Baselines.

- **Poisson Blending** [72] is implemented in OpenCV [8]. We use `cv2.seamlessClone` in the `cv2.NORMAL_CLONE` cloning type.
- **StyleGAN3**: As there is no official projection code in StyleGAN3 [45], we use unofficial implementation¹. We follow the hyperparameters of StyleGAN2 [47] official projection algorithm² and set the number of optimization iterations as 1,000.
- **StyleMapGAN** [50] introduce *stylemap*, which has spatial dimensions in the latent space. We use the official pretrained networks with 8×8 and 16×16 *stylemap* for AFHQ [21] and 32×32 *stylemap* for FFHQ [46].
- **SDEdit** [64] transforms a noise-added image into a realistic image through iterative denoising. The total denoising step is a sensitive hyperparameter that decides blending quality. If it is too small, blending results are faithful to the input images but less realistic. If it is too large, blending results are less faithful to the input images but more realistic. We carefully select the number of iterations for the best quality; 300 for the small editing parts (eyes, nose, and lip) and 500 for

the large editing parts (face and hair). To exploit the FFHQ-pretrained model [3]³, we implement guided-diffusion [30] version of SDEdit.

- **Latent Composition** [11] requires a mask to decide which area preserve. In our blending experiment, we need to preserve both the original and reference image, so we use a mask for the entire image.

Datasets. We select FFHQ [46] and AFHQv2-Cat [21] for our comparison experiments. FFHQ has 1024×1024 images, and AFHQ has 512×512 images. In SDEdit [64, 3], the pretrained model is trained on 256×256 FFHQ. Our backbone network EG3D [12] is trained on 512×512 FFHQ. In Table 1, we upsample the blending results of SDEdit and EG3D using bilinear interpolation for a fair comparison. Other baselines in FFHQ and all methods in AFHQ have the same resolution with the corresponding datasets. As the StyleMapGAN network is trained on AFHQv1, we fine-tune the networks using AFHQv2-Cat. As EG3D uses a different crop version of FFHQ compared to the original FFHQ, we fine-tune the EG3D using the original crop version of FFHQ. We use ShapeNet-Car [14] (128×128) to further demonstrate the effectiveness of our method.

Metrics. In Tables 1 and 2, we use the masked LPIPS ($LPIPS_m$) [103] to evaluate the faithfulness to the reference object. It needs a reference image to compute the score, and we use aligned reference images as pseudo-ground-truth images in both experiments: without and with 3D-aware alignment. Additionally, we always apply the alignment to our method in Tables 1 and 2.

Hyperparameters. In the 3D-aware blending, we optimize the latent code $\mathbf{w}_{\text{edit}} \in \mathbb{R}^{14 \times 512}$ of the $\mathcal{W}+$ space [1] for 200 iterations. The initial value of \mathbf{w}_{edit} is the latent code of the original image \mathbf{w}_{ori} . Adam [54] optimizer is used with 0.02 learning rate, $\beta_1 = 0.9$, and $\beta_2 = 0.999$. We reformulate Eqn. 4 to specify the hyperparameters as follows:

$$\begin{aligned} \mathcal{L}_{\text{image}} = & \|(\mathbf{1} - \mathbf{m}) \circ \mathbf{I}_{\text{edit}} - (\mathbf{1} - \mathbf{m}) \circ \mathbf{I}_{\text{ori}}\|_1 \\ & + \lambda_1 \mathcal{L}_{\text{LPIPS}}((\mathbf{1} - \mathbf{m}) \circ \mathbf{I}_{\text{edit}}, (\mathbf{1} - \mathbf{m}) \circ \mathbf{I}_{\text{ori}}) \\ & + \lambda_2 \lambda_m \mathcal{L}_{\text{LPIPS}}(\mathbf{m} \circ \mathbf{I}_{\text{edit}}, \mathbf{m} \circ \mathbf{I}_{\text{ref}}), \end{aligned} \quad (7)$$

where $\lambda_1 = 1$, $\lambda_2 = 0.5$ for AFHQ, $\lambda_2 = 0.1$ for FFHQ except hair ($\lambda_2 = 0.3$). $\lambda_m = \frac{3 \cdot H \cdot W}{\|\mathbf{m}\|_0}$ is a weighting parameter to give a high weight for the small target blending region \mathbf{m} ; H and W denotes the height and width of the image, \mathbf{m} is a binary mask $\mathbf{m} \in \{0, 1\}^{3 \times H \times W}$, and $\|\cdot\|_0$ is the L_0 norm that counts the total number of non-zero elements. Eqn. 5 is

¹<https://github.com/PDillis/stylegan3-fun>

²<https://github.com/NVlabs/stylegan2-ada-pytorch>

³<https://github.com/yandex-research/ddpm-segmentation>

reformulated as follows:

$$\begin{aligned} \mathcal{L}_{\text{density}} = & \lambda_m \sum_{\mathbf{r} \in \mathcal{R}_{\text{ref}}} \sum_{\mathbf{x} \in \mathbf{r}} \|G_{\sigma}(\mathbf{w}_{\text{edit}}; \mathbf{x}) - G_{\sigma}(\mathbf{w}_{\text{ref}}; \mathbf{x})\|_1 \\ & + \sum_{\mathbf{r} \in \mathcal{R}_{\text{ori}}} \sum_{\mathbf{x} \in \mathbf{r}} \|G_{\sigma}(\mathbf{w}_{\text{edit}}; \mathbf{x}) - G_{\sigma}(\mathbf{w}_{\text{ori}}; \mathbf{x})\|_1. \end{aligned} \quad (8)$$

The results of L1 loss in the image- and density-blending are normalized by the number of elements. We set λ in Eqn. 6 to 10.

B. Inversion details

We train our encoder to predict the camera pose of an image. It has a similar structure to the EG3D discriminator [12] but does not include minibatch discrimination [76] and involves adjustment of the number of input and output channels. Given that training images for generative NeRFs [12, 68] are pre-aligned with respect to scale and translation, the camera pose of the input image can be simplified as a rotation matrix; EG3D also samples camera poses on the surface of a sphere. As directly predicting the camera extrinsics $\in \mathbb{R}^{4 \times 4}$ is not easy, we convert the extrinsics to Euler angles $\in \mathbb{R}^3$. During inference, we transform it back to the camera extrinsics after obtaining the Euler angles from the encoder, which produces a more accurate pose estimation. Let the camera pose of the input image be \mathbf{c} .

We adopt Pivotal Tuning Inversion (PTI) [74] as our inversion method. In the first stage, we optimize the latent code $\mathbf{w} \in \mathcal{W}$ using the reconstruction loss \mathcal{L}_{rec} as follows:

$$\mathcal{L}_{\text{rec}} = \|\mathbf{I} - G_{\text{RGB}}(\mathbf{w}, \mathbf{c})\|_1 + \mathcal{L}_{\text{LPIPS}}(\mathbf{I}, G_{\text{RGB}}(\mathbf{w}, \mathbf{c})), \quad (9)$$

where \mathbf{I} is an input image and G_{RGB} is an image rendering function based on the generative NeRF G . $\|\cdot\|_1$ is the L1 distance and $\mathcal{L}_{\text{LPIPS}}$ is a learned perceptual image patch similarity (LPIPS) [103] loss. We optimize the latent code $\mathbf{w} \in \mathbb{R}^{512}$ for 300 iterations.

In the second stage, we fine-tune the generative NeRF G using the same reconstruction loss \mathcal{L}_{rec} and an additional regularization loss \mathcal{L}_{reg} as follows:

$$\begin{aligned} \mathcal{L}_{\text{reg}} = & \|G_{\text{RGB}}^f(\mathbf{w}_s, \mathbf{c}) - G_{\text{RGB}}(\mathbf{w}_s, \mathbf{c})\|_1 \\ & + \mathcal{L}_{\text{LPIPS}}(G_{\text{RGB}}^f(\mathbf{w}_s, \mathbf{c}), G_{\text{RGB}}(\mathbf{w}_s, \mathbf{c})), \end{aligned} \quad (10)$$

where G_{RGB}^f represents the frozen version of G_{RGB} , and \mathbf{w}_s is a randomly sampled latent code followed by linear interpolation with \mathbf{w} . The interpolation parameter is also randomly sampled in $[0, 1]$. The final loss function for optimizing G is as follows:

$$\mathcal{L} = \mathcal{L}_{\text{rec}} + \lambda_{\text{reg}} \mathcal{L}_{\text{reg}}, \quad (11)$$

where $\lambda_{\text{reg}} = 0.1$ and we optimize the G for 100 iterations.

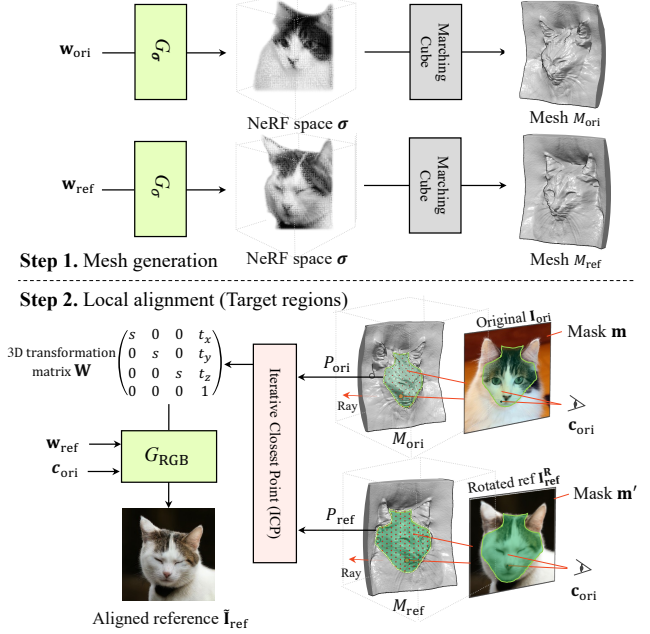


Figure 12: **Local alignment:** In **Step 1**, we first generate 3D meshes using the density field generator G_{σ} and the Marching Cube [63] algorithm. In **Step 2**, we calculate the intersected 3D points between the 2D mask \mathbf{m} and the corresponding mesh M . Then, we use Iterative Closest Point (ICP) algorithm [5, 20] to estimate the 3D transformation matrix $\mathbf{W} \in \mathbb{R}^{4 \times 4}$ to align the 3D point clouds (P_{ori} and P_{ref}) in terms of the scale and the translation. Finally, we locally align the reference image $\tilde{\mathbf{I}}_{\text{ref}}$ by $G_{\text{RGB}}(\mathbf{w}_{\text{ref}}, \mathbf{c}_{\text{ori}}; \mathbf{W})$.

Runtime. As described above, our inversion method consists of two-stage optimization. For a single image, the first stage (latent code) takes 23.6s, and the second stage (generative NeRF) takes 20.4s. We test the runtime on a single A100 GPU.

C. Local alignment

Local alignment is a fine-grained alignment between the target regions of two images. Even though we have matched two images through pose alignment, the scale and translation of target regions (*e.g.*, face, eyes, ears, *etc.*) might need to be further aligned, as the location and size of each object part differ across two object instances. Figure 12 illustrates our local alignment algorithm. In *Step 1*, we need to obtain 3D meshes M_{ori} and M_{ref} of input images using the Marching Cube algorithm [63] and the density fields. In *Step 2*, we first cast rays through the interior of the target region \mathbf{m} and determine the intersected 3D points with the mesh M_{ori} . We define the set of points as a 3D point cloud P_{ori} . Similarly, we can get the reference’s 3D point cloud P_{ref} .

Given two sets of 3D point clouds P_{ref} and P_{ori} , we use the Iterative Closest Point (ICP) algorithm [5, 20] to obtain the 3D transformation matrix $\mathbf{W} \in \mathbb{R}^{4 \times 4}$. When we sample the

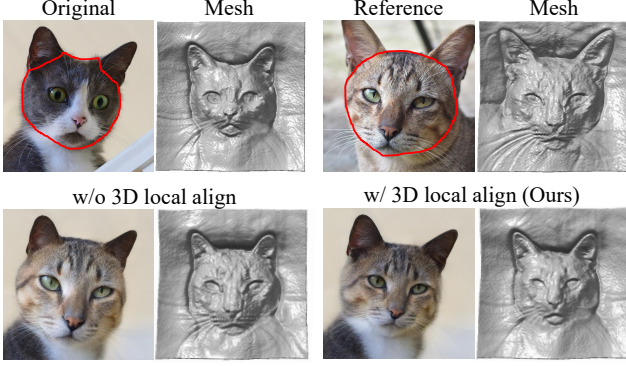


Figure 13: The effect of local alignment. Generative NeRF can align global object poses, but cannot handle differences in local parts. For example, in AFHQ-Cat, the proportions of a cat’s face and ears vary. (Bottom left) It cannot be handled with pose estimation only, so local parts can be blended in the wrong places. (Bottom right) Our 3D local alignment method alleviates this issue and produces more natural results.

3D points to align the reference image, we transform the 3D point coordinates by multiplying the coordinate of sampled points with \mathbf{W} . We only use uniform scaling and translation in ICP, as we have already matched the rotation through *pose alignment* in Section 3.1. Finally, we can generate the fully aligned reference image $\tilde{\mathbf{I}}_{\text{ref}}$ as follows:

$$\tilde{\mathbf{I}}_{\text{ref}} = G_{\text{RGB}}(\mathbf{w}_{\text{ref}}, \mathbf{c}_{\text{ori}}; \mathbf{W}). \quad (12)$$

Iterative Closest Point (ICP) [5, 20] is an algorithm to minimize the difference between two 3D point clouds. It is an iterative optimization process until we meet the threshold τ of difference or the maximum number of iterations $K = 20$. Please note that we do not use rotation in ICP as we align the rotation for entire objects in *pose alignment*, and refer to the exact details of ICP in our code.

3D transformation in the NeRF space. \mathbf{W} is the 3D transformation matrix computed by the ICP algorithm. It transforms the point cloud of the blending region in the reference, aligning with the point cloud of the blending region in the original. A Neural Radiance Field learns a mapping function f that outputs density σ and radiance \mathbf{c}_{RGB} for a given 3D point $\mathbf{x} \in \mathbb{R}^3$; $f(\mathbf{x}) = (\sigma, \mathbf{c}_{\text{RGB}})$. 3D points are sampled along a ray, which depends on the camera pose. Let us assume that we define the matrix \mathbf{W} as reducing scale. We cannot directly reduce the object scale in NeRF because the mapping function f is fixed. Instead, we can sample points more widely to reduce the object scale relatively. Renewed density σ' and radiance \mathbf{c}'_{RGB} can be computed using the inverse matrix of \mathbf{W} as follows $f(\mathbf{W}^{-1}\mathbf{x}) = (\sigma', \mathbf{c}'_{\text{RGB}})$.

Ablation study for local alignment. Pose alignment is a critical part, but it alone might not be enough to handle

loosely aligned images such as AFHQv2-Cat. Figure 13 shows an ablation study of our 3D local alignment method. The cat in the original image has a smaller face than the reference cat. If we just apply pose alignment only, the blending result (bottom left in Figure 13) will create a much bigger face than humans normally expect. After locally aligning the face of the reference, we can obtain a natural-looking result with proper scale (bottom right in Figure 13). Figure 14 shows additional ablation results. We further examine local alignment by conducting a user study. In 60.3% of cases, users prefer our blending results with local alignment compared to those without alignment. Please see the details of the user study in Appendix F.

Runtime. Pose alignment takes 1.7s, and local alignment takes 4.2s, 4.6s, 5.9s for ears, eyes, face in AFHQv2-Cat [21]. The larger blending part takes more time to align. Our local alignment implementation uses the `Trimesh`⁴ library, which provides pre-defined functions for triangular meshes on the CPU. Implementing it with GPU-based libraries can further reduce alignment runtime.

D. Details of 3D-aware Blending

Choice for the optimization space. In our blending method, we optimize the latent code $\mathbf{w} \in R^{14 \times 512}$ in the extended latent space $\mathcal{W}+$ [1]. There are other options for optimization, such as an intermediate latent space \mathcal{W} [46] or generator G . Figure 15 shows the comparison among optimization spaces in 3D-aware blending. \mathcal{W} space shows realistic blending results but lacks faithfulness to the reference object. The blended image can not capture the green eyes of the reference object, as shown in the first row of Figure 15. Optimizing G shows faithfulness to the reference object but lacks realism. The boundaries of results look particularly odd. $\mathcal{W}+$ space shows great results in both realism and faithfulness, our method adopts to optimize the $\mathcal{W}+$ latent space.

Masks used in alignment and blending. At first, we apply *pose alignment* to rotate the reference to match the pose of the original. A user selects a mask \mathbf{m} for the target blending region of the original image. The user selects another mask \mathbf{m}' for the target blending region of the reference image, or it can be derived automatically using previous works [108, 40]. We blend images using the union of the masks, $\mathbf{m} \cup \mathbf{m}'$, and we redefine \mathbf{m} as this union in Figure 4. Before blending in AFHQ, \mathbf{m} and \mathbf{m}' are used for *local alignment* to align target regions of two images, as shown in Figure 12. After local alignment, the target reference mask \mathbf{m}' is also modified.

⁴<https://trimsh.org/trimesh.html>

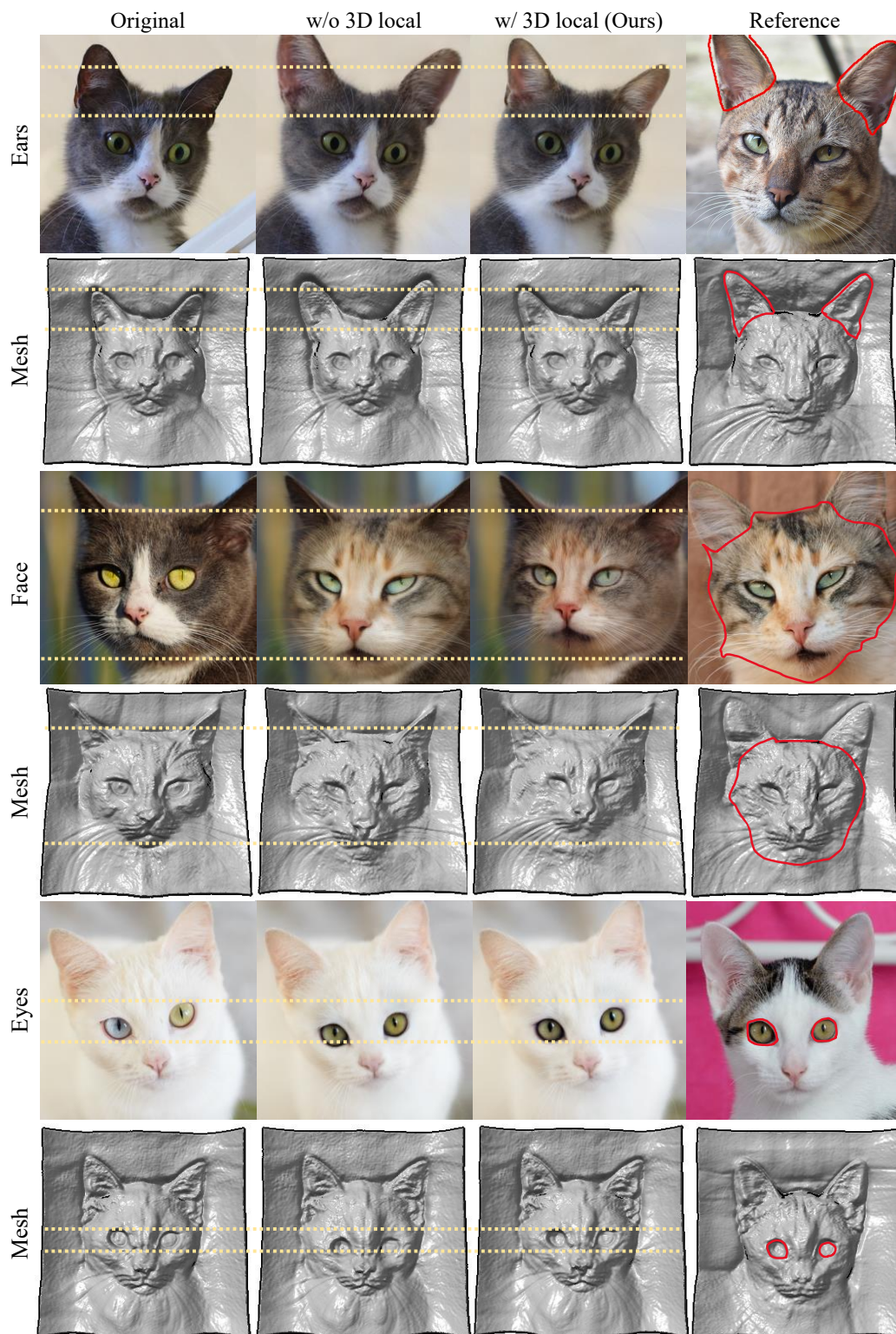


Figure 14: Ablation study of our 3D local alignment method in Appendix C. The leftmost and rightmost columns denote the original and reference images, respectively. Images in even rows are meshes corresponding to odd rows. Red lines denote the target blending mask, and yellow dotted lines are guidelines to be aware of alignment easily. Our 3D local alignment method (3rd column) shows more realistic blending results than those without local alignment (2nd column).

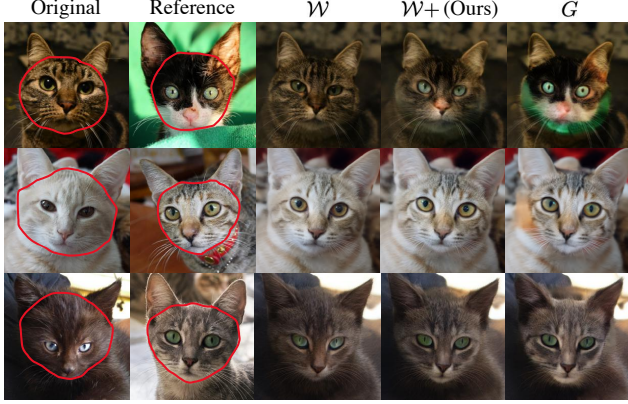


Figure 15: Blending comparison among the optimization spaces in Appendix D. The first and second columns denote original and reference images, respectively. The blending results of \mathcal{W} space (3rd column) show realistic but less faithful to the reference images; see the details of eyes. The blending results of \mathcal{G} space (rightmost column) show faithful but less realistic results. $\mathcal{W}+$ space (4th column) shows favorable results in both realism and faithfulness.

Runtime comparison with baselines. Poisson blending [72], StyleMapGAN [50], Latent Composition [11] take less than a second to blend images at 1024×1024 resolution, as they use low-level visual cues [72] or encoders [50, 11]. SDEdit [64] and StyleGAN3 [45] both require an iterative process. SDEdit takes 29.5s for 500 iterations (face and hair) and 20.0s for 300 iterations (nose, eyes, lip) at 256^2 image resolution. StyleGAN3 takes 69.8s for 512^2 and 106.5s for 1024^2 resolution. Our method takes 26.9s at 512^2 and we can further reduce the time to 12.1s by combining ours with Poisson blending (Section 3.3). All runtimes are measured in the same device with a single A100 GPU. Ours is faster than other optimization-based baselines but slower than the encoder-based methods. In the future, we can directly train an encoder using our image- and density-blending losses to reduce the runtime.

E. 3D-aware blending in StyleSDF

StyleSDF [68] generates high-fidelity view-consistent images in 1024×1024 resolution. The main difference between StyleSDF and EG3D [12] is StyleSDF uses Signed Distance Fields (SDF) as 3D representation. Besides, StyleSDF does not require camera pose labels to train the generator, unlike EG3D. In our 3D-aware blending method, we can use the SDF value d as a 3D signal, similar to the density σ in EG3D. If we assume a non-hollow surface, the SDF value can be converted to the density σ as follows:

$$\sigma(\mathbf{x}) = \frac{1}{\alpha} \cdot \text{Sigmoid}\left(\frac{-d(\mathbf{x})}{\alpha}\right), \quad (13)$$

where $\mathbf{x} \in \mathbb{R}^3$ is a 3D location and α is a learned parameter about the tightness of the density around the surface bound-



Figure 16: Our 3D-aware blending results of StyleSDF [68] in Appendix E. The third column is blending results in 1024^2 resolution. Our method can be applied to the SDF-based 3D-aware generator beyond NeRFs.

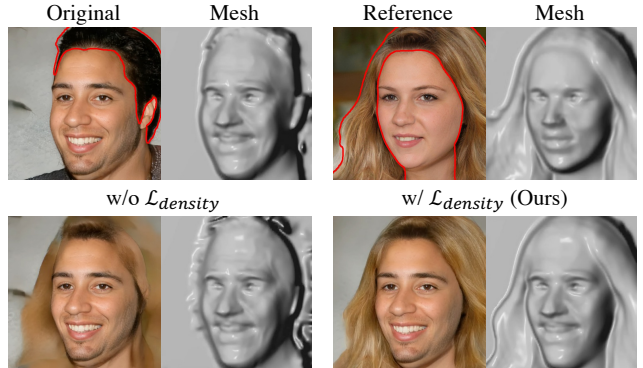


Figure 17: Ablation study of density-blending loss $\mathcal{L}_{\text{density}}$ in StyleSDF (Appendix E). This experiment uses SDF instead of volume density σ . Without $\mathcal{L}_{\text{density}}$, the blending result shows a blurry image, and the corresponding mesh can not reflect the geometry of the hair of the reference image. 3D signals such as density and SDF are key components in the 3D-aware blending method.

ary. We use the same blending loss functions used in EG3D, except we replace the density σ with SDF d .

We demonstrate that our method can be applied to other 3D-aware generative models beyond EG3D. Figure 16 shows our 3D-aware blending results in StyleSDF using generated images. Figure 17 and Figure 18 show ablation studies of our blending loss terms: $\mathcal{L}_{\text{density}}$ and $\mathcal{L}_{\text{image}}$, respectively. The ablation studies show a similar tendency with EG3D exper-

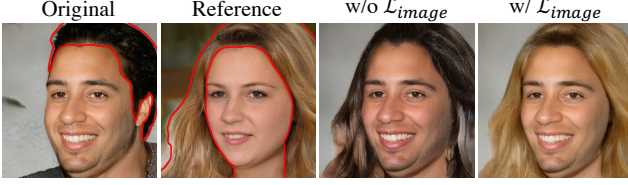


Figure 18: Ablation study of image-blending loss $\mathcal{L}_{\text{image}}$ in StyleSDF (Appendix E). We can control the degree to which the color of the reference object is reflected in the blended image, by adjusting the weight of the image-blending loss with respect to the reference image.

iments in the paper. Without the density-blending loss, we cannot blend highly structured objects like hair. If a user does not want to reflect the reference color, we remove or give a low weight to the image-blending loss on reference: λ_2 in Eqn. 4.

F. User study details

We conduct extensive user studies to show the effectiveness of our methods in the realism score of human perception on the Amazon Mechanical Turk (MTurk) platform. We refer to user study pipelines of SDEdit [64] and modify the template⁵ from the previous work [102]. The instruction page is shown in Figure 19, and MTurk workers participate in surveys comparing the results of two methods, as shown in Figures 20 and 21.

Each evaluation set consists of 25 pairwise comparison questions, with an additional five questions used to detect deceptive workers. We only invite workers with a HIT Approval Rate greater than 98%. Each set takes approximately 2 to 3 minutes and offers a reward of \$0.5.

Comparison with baselines are conducted in Tables 3 and 4. Figure 20 shows the comparison page. There are two blending results: one for our method and the other for one of the baselines. The order of images is randomly shuffled, and a worker is instructed to select a more realistic image. In CelebA-HQ [44] and AFHQv2-Cat [21], we use 60 and 40 evaluation sets, respectively; CelebA-HQ for 1,500 and AFHQ for 1,000 pairwise comparisons. We set the same number of evaluation sets to report the combination of our method and Poisson blending.

Ablation study of 3D local alignment is conducted in Appendix C. The user study page is shown in Figure 21. Experimental settings are almost similar to previous studies: Tables 3 and 4 of the main paper. We show two blending results of our method with and without *local alignment*. Note that pose alignment is used in both of the results. We use 20 evaluation sets in AFHQv2-Cat; 500 pairwise comparisons.

⁵https://github.com/phillipi/AMT_Real_vs_Fake

About this HIT:

- Please select carefully, there are some test examples to detect fake workers. If you just select one of the image randomly, you can not be approved.
- It should take about 3 minutes; 5 seconds per each comparison. There is 30 comparisons.
- You will take part in an experiment involving visual perception. You'll see a series of pairs of images. In each pair, both of images are edited images using difference AI technology.
- Your task is to determine which image looks more real.
- We only accept one hit per worker, do not do the survey twice. It will be rejected later.
- There are some unpleasant images, if you don't want to see them please turn back.

Please accept HIT to start

Figure 19: An introduction page of user studies in Amazon Mechanical Turk (MTurk).

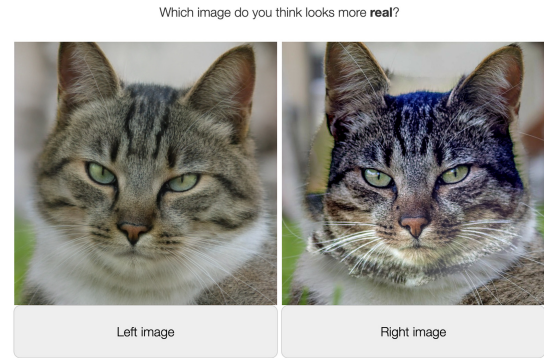


Figure 20: A comparison page between ours and baselines in MTurk.

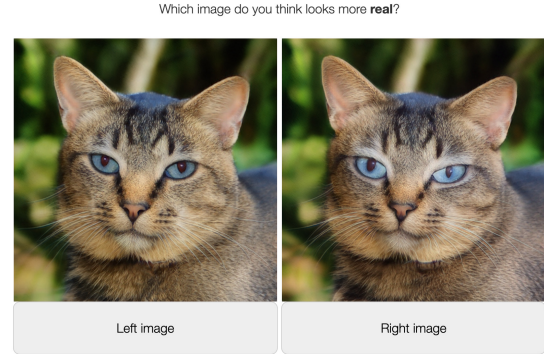


Figure 21: A comparison page of ablation study of our local alignment in MTurk.

G. Failure cases

Besides inversion as we described in the main paper, we present other failure cases in our blending method. Figure 22 shows our image blending result from a large mask to a small mask. Undesirable effects (yellow box) of the reference image have been introduced to the final result. One potential solution for future work is to inpaint the original image before blending. Figure 23 shows another failure case

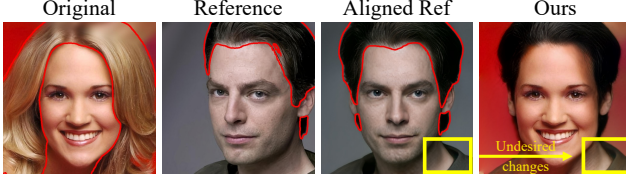


Figure 22: A failure case of blending from long hair to short hair.

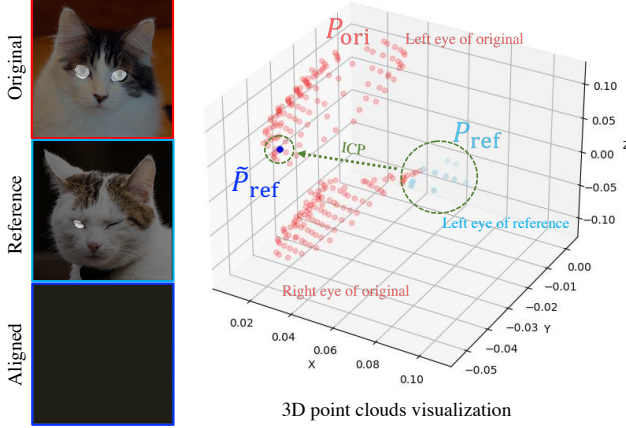


Figure 23: A failure case of the Iterative Closest Point (ICP) algorithm in Appendix G. If a user gives inappropriate masks or the target region of the mesh has indistinctive geometry, ICP may fail and generates a degraded aligned reference image. Masks are given as both eyes for the original image and the left eye for the reference image. The right figure shows the point cloud visualization. **Red** dots stand for the point cloud P_{ori} of both eyes of the original object. **Light blue** dots stand for the point cloud P_{ref} of the left eye of the reference object. **Blue** dots denote the transformed point cloud \tilde{P}_{ref} of the reference after applying ICP.

of local alignment in the Iterative Closest Point (ICP) algorithm [5, 20]. ICP is an approach to aligning the two point clouds: P_{ori} and P_{ref} for the original and reference objects, respectively. It iteratively minimizes the distance between each 3D point in P_{ori} and its nearest point in P_{ref} , but it may fall into the local extremum. P_{ref} shrinks to a point \tilde{P}_{ref} as discussed in the previous work [90]. To mitigate this issue, we restrict the minimum and maximum scaling factor to [0.75, 1.25]. In the future, we may utilize recent pairwise registration techniques [22, 90] instead of ICP for better local alignment.

H. Societal impact

The societal impact of image blending shares similar issues raised in the other generative models and view synthesis techniques [38]. For instance, a malicious user may use image blending to manipulate the expression and identity of a real person or create a scene that does not exist in real life. Most of the existing watermarking and visual forensics methods focus on 2D image content [97]. Adopting visual

forensics methods for 3D-aware content seems to be an important future work.

I. Additional results

Comparison with StyleFusion and Barbershop. In order to show the advantages of our 3D-aware approach, we compare our method with additional baselines: StyleFusion [43] and Barbershop [109]. Previous methods struggle to blend fairly misaligned images, and they acknowledged this limitation in their papers. As shown in Figure 24, latent-based methods often fail to preserve identity, as projecting images into low-dimensional latent space remains challenging. For example, Latent Composition [11] and StyleFusion altered identities from the input images. Other works, such as StyleMapGAN [50] and Barbershop, use the spatial latent space to capture image details, but as a trade-off, they are worse at blending misaligned images. In contrast, our method preserves identity while maintaining 3D consistency between parts. Our method addresses these issues by 1) 3D-aware alignment and 2) blending with 3D-aware constraints, including pixel RGBs and volume density from the aligned reference.

Comparison with SDEdit in AFHQ. In addition to comparing CelebA-HQ [44], we also conduct a comparison with SDEdit using the AFHQv2-Cat dataset [21]. However, since pretrained diffusion models on AFHQ are not publicly available, we use the LSUN-Cat model⁶. Figure 25 shows SDEdit fails to preserve the reference well in both datasets.

More qualitative results. We show additional qualitative results of our 3D-aware blending. Figure 26 shows blending results of highly structured objects such as eyeglasses. Figures 27 and 28 show blending comparisons with baselines. Our results show outstanding blending results than other baselines. By virtue of NeRF, we can synthesize novel-view images. Figures 29–32 show multi-view consistent blending results. Please see our project page for the multi-view consistent blending videos.

⁶<https://github.com/openai/guided-diffusion>

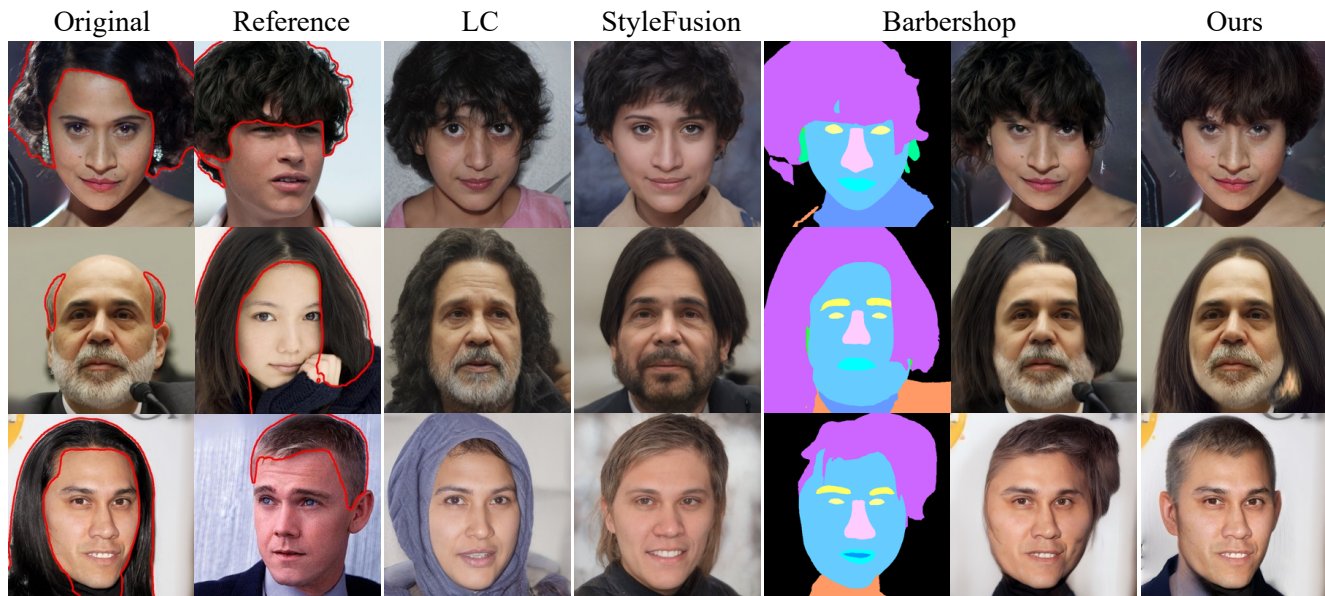


Figure 24: Additional hair blending comparison with baselines. Latent Composition (LC) generates real-looking images but far from the input images. StyleFusion alters the hair length of the reference images. In StyleFusion and Barbershop, the hair and face poses do not match in each blending result. Ours shows the best results regarding 3D-aware alignment and identity preservation.

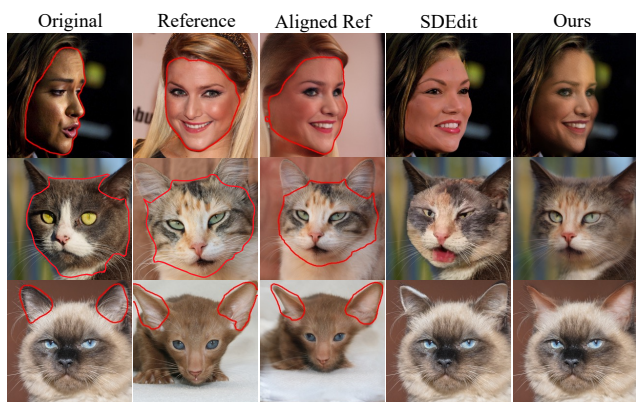


Figure 25: Comparison with SDEdit in CelebA-HQ and AFHQv2-Cat test sets. SDEdit blending results often do not match reference images well. The identities of faces or ears have been changed in both CelebA-HQ and AFHQ datasets.

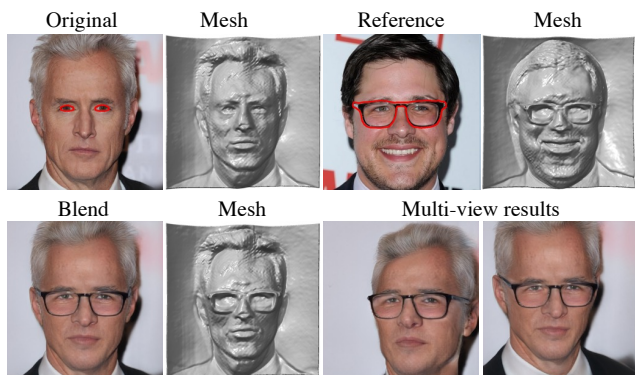


Figure 26: Blending highly structured parts such as eyeglasses. Our blending method can reflect the high-fidelity 3D shape of eyeglasses and generate multi-view consistent results.



Figure 27: Blending comparison with baselines in the CelebA-HQ test set. It shows supplementary comparison results of Table 1. LC and PB stand for Latent Composition [11] and Poisson Blending [72], respectively.



Figure 28: Blending comparison with baselines in the AFHQv2-Cat test set. It shows supplementary comparison results of Table 2.

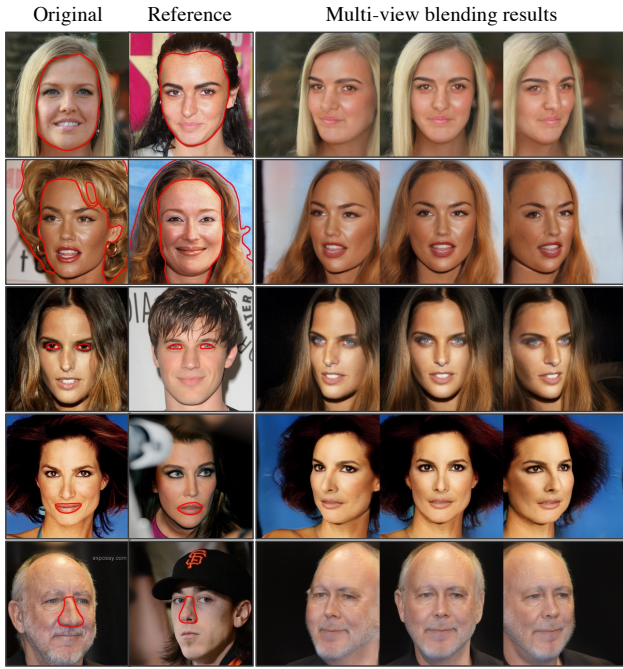


Figure 29: Multi-view blending results in CelebA-HQ using EG3D.

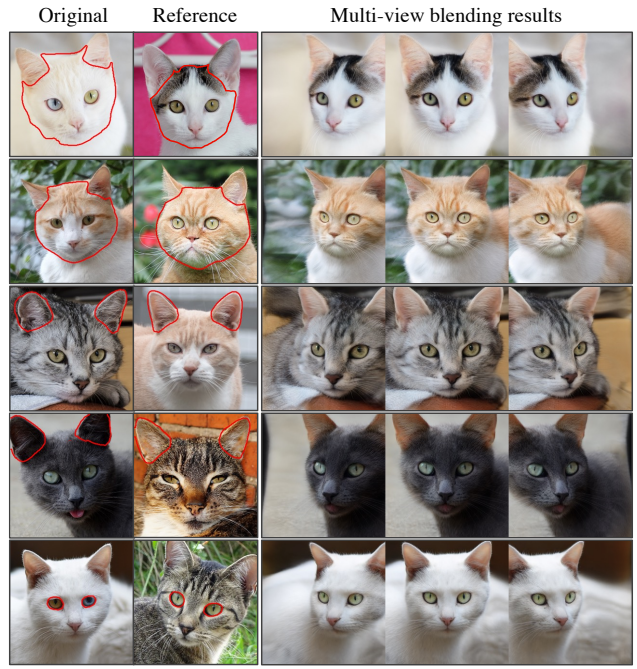


Figure 31: Multi-view blending results in AFHQv2-Cat using EG3D.



Figure 30: Multi-view blending results in ShapeNet-Car [14] using EG3D. The last two rows demonstrate that our method can achieve natural blending even when the sizes of the blending regions differ. Our 3D local alignment also performs well on the ShapeNet-Car dataset.

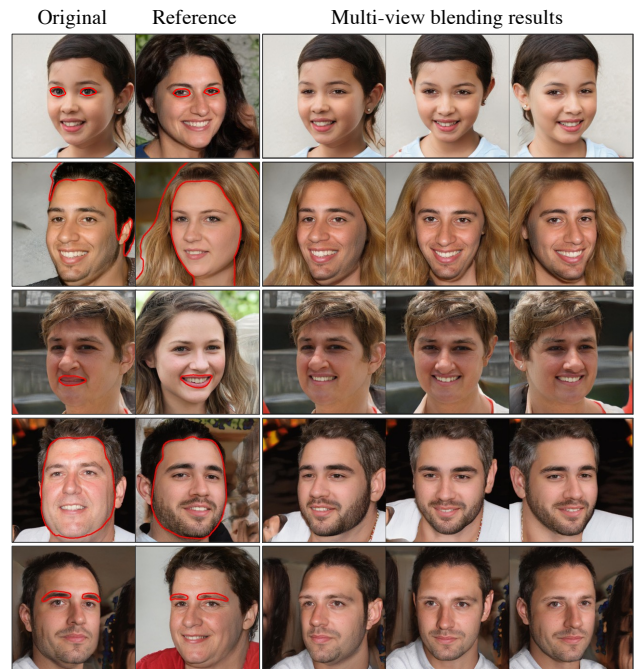


Figure 32: Multi-view blending results in generated images using FFHQ-pretrained StyleSDF [68].