

---

# ZIGNeRF: Zero-shot 3D Scene Representation with Invertible Generative Neural Radiance Fields

---

Kanghyeok Ko                                  Minhyeok Lee<sup>\*</sup>  
 dogworld12@cau.ac.kr                                  mlee@cau.ac.kr

School of Electrical and Electronics Engineering, Chung-Ang University

## ABSTRACT

Generative Neural Radiance Fields (NeRFs) have demonstrated remarkable proficiency in synthesizing multi-view images by learning the distribution of a set of unposed images. Despite the aptitude of existing generative NeRFs in generating 3D-consistent high-quality random samples within data distribution, the creation of a 3D representation of a singular input image remains a formidable challenge. In this manuscript, we introduce ZIGNeRF, an innovative model that executes zero-shot Generative Adversarial Network (GAN) inversion for the generation of multi-view images from a single out-of-domain image. The model is underpinned by a novel inverter that maps out-of-domain images into the latent code of the generator manifold. Notably, ZIGNeRF is capable of disentangling the object from the background and executing 3D operations such as 360-degree rotation or depth and horizontal translation. The efficacy of our model is validated using multiple real-image datasets: Cats, AFHQ, CelebA, CelebA-HQ, and CompCars.

## 1 Introduction

The remarkable success of generative adversarial networks (GANs) [8] has spurred significant advancements in realistic image generation with high quality. Particularly, following the emergence of StyleGAN [17], numerous 2D-based generative adversarial network models have benefited from a deeper understanding of latent spaces [16, 18]. Consequently, various computer vision tasks, such as conditional image generation and style transfer [13, 20], have shown substantial progress. However, 2D-based image generation models are constrained in their ability to generate novel view images due to their limited understanding of the underlying 3D geometry of real-world scenes.

To overcome this challenge, several studies have adopted the neural radiance field (NeRF) [24] approach, which encodes a scene into a multi-layer perceptron (MLP) to provide 3D rendering. Although conventional NeRF [24] has successfully facilitated the development of 3D-aware models and reduced computational costs in novel view synthesis tasks, it remains impractical to train a model overfitted to a single scene with multi-view images [24, 47]. Consequently, various studies have extended NeRF by integrating it with generative models, i.e., generative NeRF. Generative NeRF [1, 2, 6, 9, 27, 28, 34] models can be trained on unposed real-world images, whereas conventional NeRF necessitates multiple images of a single scene [37, 39, 44]. Moreover, generative NeRF has been employed for obtaining conditional samples through techniques such as class label information [14] or text encoding [7, 29, 30, 41].

Despite the convenience and intuitiveness of these approaches, they possess limitations in image editing and generating 3D representations of specific inputs, such as out-of-domain images or real-world images. To enable more practical applications, generative NeRF models have also

---

<sup>\*</sup>Corresponding author.

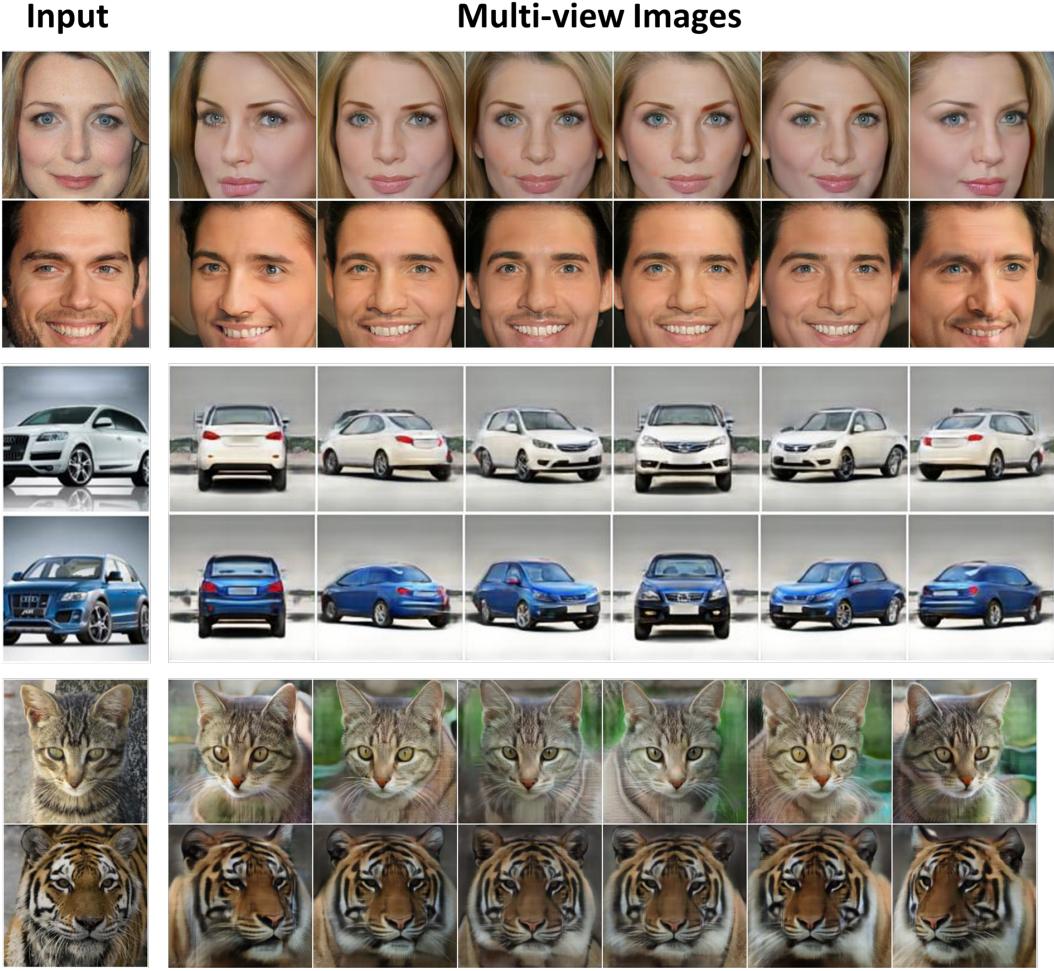


Figure 1: **Demonstration of the 3D reconstruction results employing our proposed method, ZIGNeRF.** This illustration depicts the successful zero-shot 3D GAN inversion across various real-world image datasets [5, 15, 45].

incorporated GAN inversion techniques [31, 32, 38, 50] for the 3D representation of particular input images, including out-of-distribution or real-world images. However, previous studies have faced a constraint that necessitates fine-tuning on pre-trained models for specific images [19, 21, 42, 46]. This requirement hinders the application of these models to numerous real samples simultaneously and renders the process time-inefficient, as it demands extensive fine-tuning.

In this study, we propose a novel zero-shot methodology for the generation of multi-view images, derived from input images unseen during the training process. This approach leverages a 3D-aware GAN inversion technique. Notably, our model proffers 3D-consistent renderings of unposed real images during inference, eliminating the need for supplementary fine-tuning.

Our architectural design bifurcates into two distinct components: the 3D-generation module and the 3D-aware GAN inversion module. The former is founded on the principles of GIRAFFE [27], which successfully amalgamates the compositional attributes of 3D real-world scenes into a generative framework. To enhance the precision of 3D real-world reconstruction and improve image quality, we introduce modifications to the GIRAFFE module, specifically in the decoder and neural renderer. The 3D-aware GAN inverter, on the other hand, is trained with images synthesized from the generator. This strategic approach enables the inverter to accurately map the input image onto the generator's manifold, regardless of the objects' pose. Example results of our model is displayed in Fig. 1.

We subject our model to rigorous evaluation, utilizing five diverse datasets: Cats, CelebA, CelebA-HQ, AFHQ, and CompCars. Additionally, we demonstrate the model’s robustness by inputting FFHQ images into a model trained on CelebA-HQ. The primary contributions of this work are as follows:

- We present ZIGNeRF, a pioneering approach that delivers a 3D-consistent representation of real-world images via zero-shot estimation of latent codes. To our knowledge, this is the first instance of such an approach in the field.
- ZIGNeRF exhibits robust 3D feature extraction capabilities and remarkable controllability with respect to input images. Our model can perform 3D operations, such as a full 360-degree rotation of real-world car images, a feat not fully achieved by many existing generative NeRF models.

## 2 Related Work

### 2.1 Neural Radiance Field (NeRF)

NeRF is an influential method for synthesizing photorealistic 3D scenes from 2D images. It represents a 3D scene as a continuous function using a multi-layer perceptron (MLP) that maps spatial coordinates to RGB and density values, and then generates novel view images through conventional volume rendering techniques. Consequently, NeRF significantly reduces computational costs compared to existing voxel-based 3D scene representation models [11, 26, 35, 37, 49]. However, the training method of NeRF, which overfits a single model to a single scene, considerably restricts its applicability and necessitates multiple structured training images, including camera viewpoints [3, 37].

### 2.2 Generative NeRF

Generative NeRFs optimize networks to learn the mapping from latent code to 3D scene representation, given a set of unposed 2D image collections rather than using multi-view supervised images with ground truth camera poses. Early attempts, such as GRAF [34] and pi-GAN [2], demonstrated promising results and established the foundation for further research in the generative NeRF domain. Recent works on generative NeRF have concentrated on generating high-resolution 3D-consistent images. The recently proposed StyleNeRF [9] successfully generates high-resolution images by integrating NeRF into a style-based generator, while EG3D [1] exhibits impressive results with a hybrid architecture that improves computational efficiency and image quality.

However, real-life applications frequently necessitate conditional samples that exhibit the desired attribute rather than random samples in data distribution. We adopt GAN inversion as a conditional method, as opposed to class-based or text encoding conditional methods, which are prevalent in 2D generative models [4]. The aforementioned conditional generation techniques, such as class-based or text encoding methods, possess limitations. Firstly, the training dataset must include conditional information, such as labels or text corresponding to each sample. Secondly, they cannot provide 3D representation of real-world images as conditional input. We address these limitations in existing conditional generative NeRF models by introducing GAN Inversion into generative NeRF for conditional generation.

### 2.3 3D aware GAN inversion

With the remarkable progress of GANs, numerous studies have endeavoured to understand and explore their latent space to manipulate the latent code meaningfully. GAN inversion represents the inverse process of the generator in GANs. Its primary objective is to obtain the latent code by mapping a given image to the generator’s latent space. Ideally, the latent code optimized with GAN inversion can accurately reconstruct an image generated from the pre-trained generator. The output sample can be manipulated by exploring meaningful directions in the latent space [36]. Moreover, real-world images can be manipulated in the latent space using GAN inversion.

Several studies have investigated 3D GAN inversion with generative NeRF to generate multi-view images of input samples and edit the samples in 3D manifolds. Most previous works fine-tuned the pre-trained generator due to the utilization of optimization-based GAN inversion methods. However,

▪ Training process of the 3D generation part

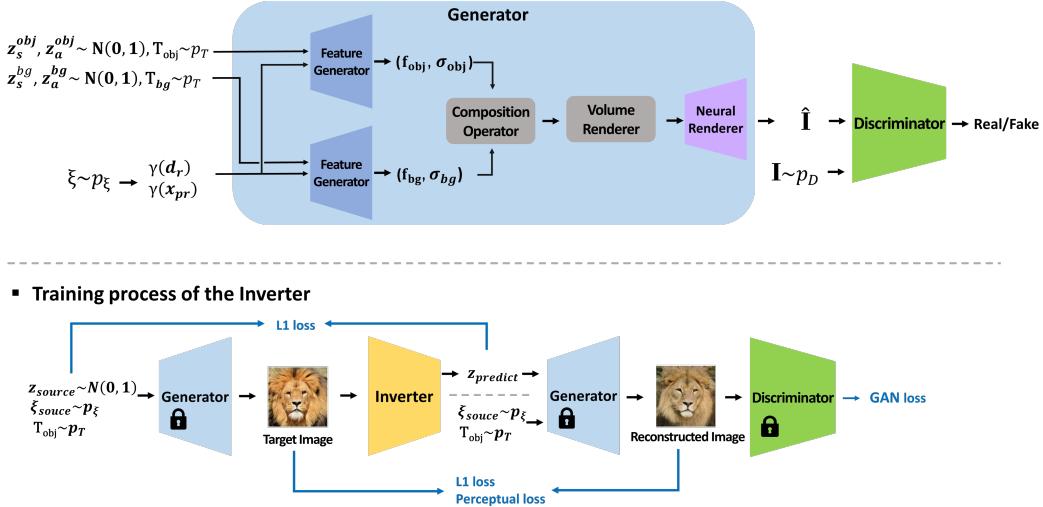


Figure 2: **The comprehensive architecture of ZIGNeRF.** The 3D generative component is trained to produce photorealistic images consistent with 3D structures by mapping the latent code and camera pose to a synthetic image. Subsequently, the inverter is trained in conjunction with the pre-trained generator and discriminator.

additional steps for fine-tuning the generator for GAN inversion impose limitations in terms of adaptability and computational costs.

In this paper, we propose a novel inverter for 3D-aware zero-shot GAN inversion. The proposed inverter can map out-of-domain images into the latent space of the generator. Our model can generate 3D representations of real-world images without requiring additional training steps. The proposed 3D-aware zero-shot GAN inversion maximizes applicability since the trained model can be directly applied to out-of-domain images.

### 3 Method

This work seeks to generate multi-view images from an out-of-domain image by combining generative NeRF with GAN inversion. The proposed method, graphically delineated in Fig. 2, encompasses two distinct phases: the 3D-generation segment and the 3D-aware inverter. The first phase involves training the 3D-generation component, an architecture based on GIRAFFE, augmented by enhancements in the neural renderer and the discriminator modules to fortify and expedite the training process. In the second phase, the 3D-aware inverter is trained with the pre-trained generator. The novel inverter is designed to transform out-of-domain images into latent codes within the generator’s latent space. Consequently, the generator can produce multi-view images of the out-of-domain image using the latent code derived from the inverter. Throughout the training of the inverter, we utilize the images generated from the generator, imbued with 3D information, as the training dataset. At test time, the inverter executes zero-shot inversion on real-world images, obviating the need for additional fine-tuning for unseen images. The proposed method thereby holds great promise for generating 3D-consistent multi-view images from real-world input images.

#### 3.1 3D Generation

**Compositional Generative Neural Feature Field.** Our 3D-generator represents a scene with a compositional generative neural feature field, a continuous function inherited from GIRAFFE, to represent a scene. This is essentially a combination of feature fields, each representing an object in a single scene, with the background also considered an object. In the 3D-generator, a 3D location,

$\mathbf{x} \in \mathbb{R}^3$ , a viewing direction,  $\mathbf{d} \in \mathbb{S}^2$ , and latent code,  $\mathbf{z} \sim \mathcal{N}(0, 1)$ , are mapped to a volume density  $\sigma \in \mathbb{R}^+$  and a high-dimensional feature field  $\mathbf{f} \in \mathbb{R}^{M_f}$ , rather than RGB colour  $\mathbf{c} \in \mathbb{R}^3$ .

Affine transformation is applied to objects in the scene so that each object can be controlled in terms of poses, which include scale, translation, and rotation:

$$T = \{\mathbf{s}, \mathbf{t}, \mathbf{R}\}, \quad (1)$$

where  $\mathbf{s}, \mathbf{t} \in \mathbb{S}$  indicate scale and translation parameters, respectively, and  $\mathbf{R} \in \text{SO}(3)$  determine rotation. The affine transformation enables object-level control by generating the bounding box corresponding to  $T$  of a single object:

$$\tau = \mathbf{R} \cdot \mathbf{s} \mathbf{I} \cdot + \mathbf{t}, \quad (2)$$

where  $\mathbf{I}$  is the  $3 \times 3$  identity matrix. Compositional generative neural feature field is parameterized with an MLP as follows:

$$C((\sigma_i, \mathbf{f}_i)_{i=1}^N) = C(f_{\theta i}(\gamma(\tau^{-1}(\mathbf{x})), \gamma(\tau^{-1}(\mathbf{d})), \mathbf{z}_i)_{i=1}^N), \quad (3)$$

$$\mathbf{z} = [\mathbf{z}_s^1, \mathbf{z}_a^1, \dots, \mathbf{z}_s^N, \mathbf{z}_a^N], \quad (4)$$

where  $\gamma(\cdot)$  is positional encoding function [24], which is applied separately to  $\mathbf{x}$  and  $\mathbf{d}$ , and  $C(\cdot)$  is the compositional operator that composites feature field from the  $N-1$  objects and a background. We then volume render the composited volume density and feature field rather than directly output the final image. 2D-feature map, which is fed into neural renderer for final synthesized output, is attained by volume rendering function  $\pi_v$ ,

$$\pi_v(C(\sigma, \mathbf{f})) = \mathbf{F}. \quad (5)$$

**Neural renderer with residual networks.** Our model outputs final synthetic image with neural rendering on the output feature map of volume rendering. We observe that the original neural renderer of GIRAFFE does not preserve the feature well. Furthermore, the learning rate of the decoder and the neural renderer is not synchronized; hence the training of the generator is unstable.

We improve the simple and unstable neural renderer of GIRAFFE. Our neural renderer replaces  $3 \times 3$  convolution layer blocks with residual blocks [10] and employs the ReLU activation rather than leaky ReLU activation [43] for faster and more effective rendering. To stabilize the neural rendering, we adopt spectral normalization [25] as weight normalization. We experimentally verify that the modified neural renderer improves the stability of the training and the quality of the outputs. Our neural renderer, which maps the feature map  $\mathbf{F}$  to the final image  $\hat{\mathbf{I}} \in \mathbb{R}^{H \times W \times 3}$ , is parameterized as:

$$\pi_\theta(\mathbf{F}) = \hat{\mathbf{I}}. \quad (6)$$

**Discriminator.** As the vanilla GAN [8], the discriminator outputs probability, which indicates whether the input image is real or fake. We replace the 2D CNN-based discriminator with residual blocks employing spectral normalization as weight normalization.

**Objectives.** The overall objective function of the 3D-generative part is:

$$L_{G, D} = L_{\text{GAN}} + \lambda L_{\text{RI}}, \quad (7)$$

where  $\lambda = 10$ . We use GAN objective [8] with R1 gradient penalty [23] to optimize the network.

### 3.2 3D-aware Invertor

To invert a given image into latent codes within the generator’s latent space, we introduce a novel inverter. This inverter is designed by stacking the residual encoder block with ReLU activations, as depicted in Fig. 3. Four linear output layers are situated at the culmination of the inverter to facilitate output. These residual blocks extract the feature of the input image, and each linear output layer estimates the  $\mathbf{z}_s^{\text{obj}}, \mathbf{z}_a^{\text{obj}}, \mathbf{z}_s^{\text{bg}}, \mathbf{z}_a^{\text{bg}}$  of the input image.

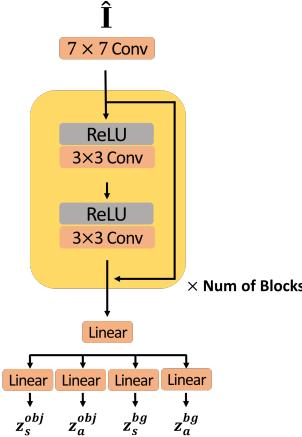


Figure 3: **Schematic representation of the architecture of the inverter deployed in ZIGNeRF.**

The challenge of 3D-aware GAN inversion involves mapping multi-view images of a single object into a unique latent code. To construct a 3D-aware inverter, we opt to use the synthesized image  $\hat{\mathbf{I}}$  as the training data. Given that we already possess the source parameters of the generated image, the inverter solely estimates the latent code  $\mathbf{z}^{\text{predict}}$  of the input image. The generated training images equip the inverter to extract the feature of unseen images, which vary in viewing direction, scale, and rotation. Following the latent code inference, the pre-trained generator reconstructs the input image using  $\mathbf{z}^{\text{predict}}$  and source parameters, which include camera pose,  $\xi^{\text{source}}$ , and compositional parameter,  $\mathbf{T}^{\text{source}} = \{\mathbf{s}, \mathbf{t}, \mathbf{R}\}$ :

$$I_{\theta}(\hat{\mathbf{I}}) = \mathbf{z}^{\text{predict}}, \quad (8)$$

$$G_{\theta}(\mathbf{z}^{\text{predict}}, \mathbf{T}^{\text{source}}, \xi^{\text{source}}) = \hat{\mathbf{I}}^{\text{reconst}}. \quad (9)$$

As the inverter learns to estimate the latent source code, we found that the L1 loss between the two latent codes in latent space was inadequate for reconstructing the scene. Thus, we opted to employ GAN loss and L1 as an image-level loss to generate a plausible image. In addition, we incorporated two perceptual losses, namely the Structural Similarity Index Measure (SSIM) [40] and the Learned Perceptual Image Patch (LPIPS) [48] loss, to conserve the fine details of the source image. The inverter can be optimized using the following function:

$$\begin{aligned} L_I = & L_{\text{GAN}}(\hat{\mathbf{I}}^{\text{predict}}) + \lambda_1 L_{\text{latent}}(\mathbf{z}^{\text{source}}, \mathbf{z}^{\text{predict}}) \\ & + \lambda_2 L_{\text{reconst}}(\hat{\mathbf{I}}^{\text{source}}, \hat{\mathbf{I}}^{\text{predict}}) + \lambda_3 L_{\text{percept}}(\hat{\mathbf{I}}^{\text{source}}, \hat{\mathbf{I}}^{\text{predict}}), \end{aligned}$$

where  $\hat{\mathbf{I}}^{\text{predict}}$  indicates the image reconstructed by the pre-trained generator using  $\mathbf{z}^{\text{predict}}$ .  $L_{\text{latent}}$  and  $L_{\text{reconst}}$  represent latent-level and image-level loss, respectively, both utilizing L1 loss.  $L_{\text{percept}}$  signifies image-level perceptual loss, employing the LPIPS loss and SSIM loss.

### 3.3 Training specifications

During the training phase, we randomly sample the latent codes  $\mathbf{z}_s, \mathbf{z}_a \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$ , and a camera pose  $\xi \sim p_{\xi}$ . The parameters  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$  are set to 10, 100, and 1, respectively, for training the inverter. The model is optimized using the RMSProp optimizer [33], with learning rates of  $1 \times 10^{-4}$ ,  $7 \times 10^{-5}$ , and  $1 \times 10^{-4}$  for the generator, the discriminator, and the inverter, respectively. We utilize a batch size of 32. For the first 100,000 iterations, the generator and the discriminator are trained, and the inverter is trained for the next 50,000 iterations. During the training process of the inverter, the generator and the discriminator remain frozen.

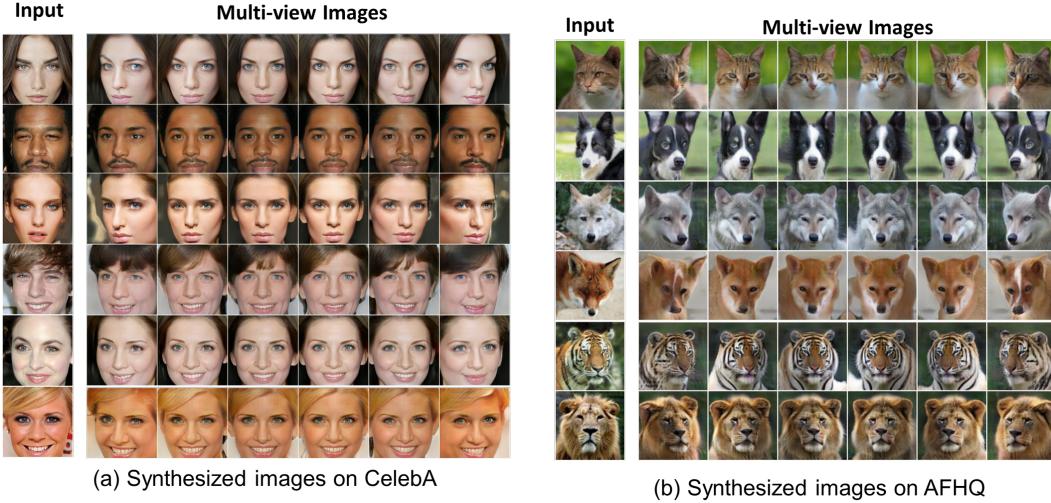


Figure 4: **Display of  $256^2$  multi-view synthesis applied to facial datasets: CelebA-HQ [15] and AFHQ [5].**



Figure 5: **Visualisation of reconstructed images based on an input car image [45], following compositional operations.** These illustrations highlight the effective disentanglement of the object from the background and the provision of 3D controllability.

## 4 Experiments

ZIGNeRF is evaluated concerning zero-shot feature extraction, 3D controllability, and adaptability. We test on five real-world datasets: Cats, AFHQ [5], CelebA [22], CelebA-HQ [15], and CompCar [45]. An additional dataset, FFHQ [17], is used to demonstrate the robust adaptation capabilities of the proposed model. All input images shown in this section were not used during the training process, thereby validating the zero-shot 3D GAN inversion with unseen images. We commence with a visual validation of the proposed model, examining the similarity between the input image and the reconstructed images and 3D-consistent controllability. The model is then evaluated using Fréchet Inception Distance (FID) [12] as a metric. We conclude with ablation studies to validate the efficacy of the loss function in optimizing the inverter.

### 4.1 Controllable 3D Scene Reconstruction

We visually demonstrate that our proposed model generates multi-view consistent images corresponding to the input image. Fig. 4 showcases 3D reconstruction on CelebA-HQ [15] and AFHQ [5], substantiating that the inverter successfully extracts facial features irrespective of gender or skin colour in human faces, and species in animal faces. Fig. 5 exhibits the model’s controllability and object disentanglement with CompCar [45], indicating that the inverter estimates the latent code of the object and background effectively. Notably, the proposed model can facilitate 3D-consistent 360-degree rotation, a common limitation of generative NeRF methods. We further attest to the robustness of our model by applying it to FFHQ, as shown in Fig. 6.

| Method        | Models        | Cats             |                  | CelebA(HQ)       |                  | CompCar          |                  | AFHQ             |                  |
|---------------|---------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|------------------|
|               |               | 128 <sup>2</sup> | 256 <sup>2</sup> |
| Unconditional | GIRAFFE       | 24.01            | 21.28            | 19.45            | 23.14            | 38.91            | 40.84            | 35.03            | 38.18            |
|               | ZIGNeRF(ours) | <b>12.31</b>     | <b>11.21</b>     | <b>11.01</b>     | <b>14.98</b>     | <b>22.67</b>     | <b>22.57</b>     | <b>12.81</b>     | <b>19.96</b>     |
| Conditional   | ZIGNeRF(ours) | <b>15.06</b>     | <b>16.83</b>     | <b>14.77</b>     | <b>25.66</b>     | <b>25.97</b>     | <b>25.41</b>     | <b>14.02</b>     | <b>28.78</b>     |

Table 1: **Comparative analysis of the FID between our proposed ZIGNeRF and a baseline model.** The models were trained on four distinct datasets with the resolution of 128<sup>2</sup> and 256<sup>2</sup>.

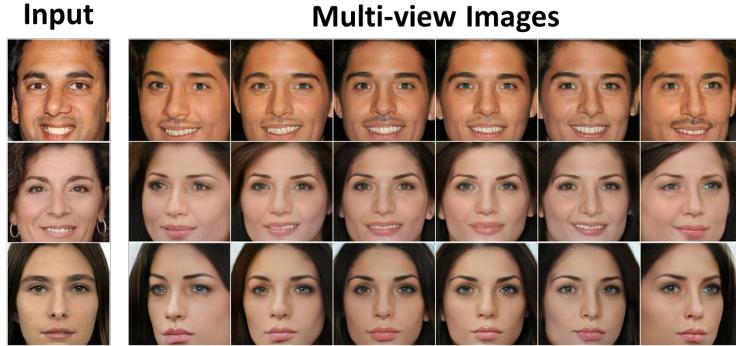


Figure 6: **Presentation of 256<sup>2</sup> synthesized images conditioned on input FFHQ [17] images, produced by the model trained on the CelebA-HQ dataset [15].**

## 4.2 Quantitative Evaluation

To thoroughly evaluate the efficacy of our proposed model, ZIGNeRF, we conduct experiments in both conditional and unconditional generation modes. The evaluation process involves a random sampling of 20,000 real images alongside 20,000 synthesized images, which is a conventional method to compare generative models. The results are displayed in Tab. 1.

In the context of the unconditional model, we generate samples using random latent codes. The training process entails 100,000 iterations. Notably, our model, ZIGNeRF, significantly outperforms the baseline GIRAFFE [27] model. As an illustration, for the CelebA(HQ) 256<sup>2</sup> dataset, ZIGNeRF achieves a score of 14.98, substantially lower than the GIRAFFE’s score of 23.14. This is indicative of the model’s ability to produce higher-quality images with fewer iterations.

Turning to the conditional synthesis, the latent codes estimated by the inverter are employed on randomly sampled real images. The training process for the generator is conducted over 100,000 iterations, while the inverter training comprises 50,000 iterations, during which the generator is kept static. When compared to GIRAFFE, ZIGNeRF demonstrate superior performance in conditional samples as well. For instance, in the AFHQ 128<sup>2</sup> dataset, our model attains a score of 14.02, marking a significant improvement over the GIRAFFE’s score of 35.03.

## 4.3 Ablation study

In the interest of validating the loss function deployed in training the inverter, we undertake an ablation study. The study scrutinizes the necessity of each loss component: latent loss, reconstruction loss, GAN loss, and perceptual loss. The imperative nature of each loss function is demonstrated through its incremental addition to the naive model, which is trained solely via latent code comparison. Fig. 7 illustrates the individual contribution of each loss function. It is observed that the naive model exhibits limited capability in reconstructing the input image. The reconstruction loss  $L_{\text{reconst}}$  aligns the reconstructed image with the input at an image-level. The GAN loss  $L_{\text{GAN}}$  is observed to enhance the realism of the reconstructed image, independent of improving the input-reconstructed image

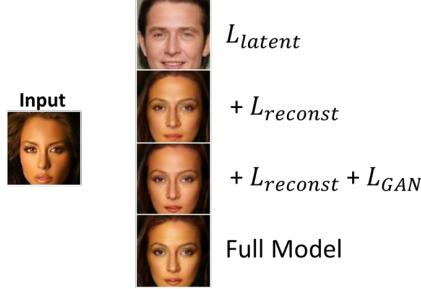


Figure 7: **Ablation study of the loss functions employed in the training of the inverter within ZIGNeRF.**

similarity. The full model elucidates that the perceptual loss  $L_{\text{percept}}$  plays a pivotal role in refining the expression of minute attributes, skin colour, and texture.

## 5 Conclusion

In this paper, we have proposed ZIGNeRF, an innovative technique that manifests a 3D representation of real-world images by infusing a 3D-aware zero-shot GAN inversion into generative NeRF. Our inverter is meticulously designed to map an input image onto a latent manifold, a learning process undertaken by the generator. During testing, our model generates a 3D reconstructed scene from a 2D real-world image, employing a latent code ascertained from the inverter. Rigorous experiments conducted with four distinct datasets substantiate that the inverter adeptly extracts features of input images with varying poses, thereby verifying the 3D controllability and immediate adaptation capabilities of our model.

Our novel approach carries the potential for wide application, given that our pipeline can be generally applied to other existing generative NeRFs. It is worth noting that this zero-shot approach is a pioneering contribution to the field, bringing forth a paradigm shift in 3D image representation. In future work, we envisage extending the proposed method by manipulating the inverted latent code for editing the input image, thereby further enhancing the capabilities of this innovative model.

## References

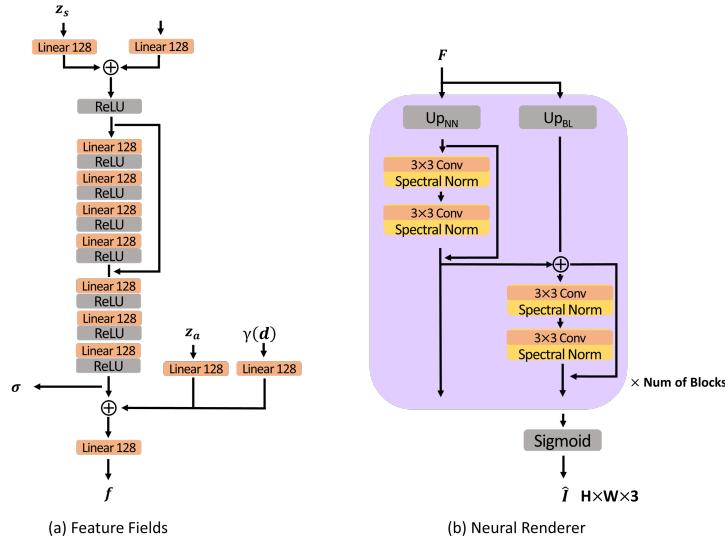
- [1] E. R. Chan, C. Z. Lin, M. A. Chan, K. Nagano, B. Pan, S. De Mello, O. Gallo, L. J. Guibas, J. Tremblay, S. Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16123–16133, 2022.
- [2] E. R. Chan, M. Monteiro, P. Kellnhofer, J. Wu, and G. Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5799–5809, 2021.
- [3] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015.
- [4] Y. Choi, M. Choi, M. Kim, J.-W. Ha, S. Kim, and J. Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8789–8797, 2018.
- [5] Y. Choi, Y. Uh, J. Yoo, and J.-W. Ha. Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8188–8197, 2020.
- [6] Y. Deng, J. Yang, J. Xiang, and X. Tong. Gram: Generative radiance manifolds for 3d-aware image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10673–10683, 2022.
- [7] R. Gal, O. Patashnik, H. Maron, A. H. Bermano, G. Chechik, and D. Cohen-Or. Stylegan-nada: Clip-guided domain adaptation of image generators. *ACM Transactions on Graphics (TOG)*, 41(4):1–13, 2022.
- [8] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- [9] J. Gu, L. Liu, P. Wang, and C. Theobalt. Stylenet: A style-based 3d-aware generator for high-resolution image synthesis. *arXiv preprint arXiv:2110.08985*, 2021.
- [10] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [11] P. Henzler, N. J. Mitra, and T. Ritschel. Escaping plato’s cave: 3d shape from adversarial rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9984–9993, 2019.
- [12] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- [13] L. Höllerin, J. Johnson, and M. Nießner. Stylemesh: Style transfer for indoor 3d scene reconstructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6198–6208, 2022.
- [14] K. Jo, G. Shim, S. Jung, S. Yang, and J. Choo. Cg-nerf: Conditional generative neural radiance fields. *arXiv preprint arXiv:2112.03517*, 2021.
- [15] T. Karras, T. Aila, S. Laine, and J. Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017.
- [16] T. Karras, M. Aittala, S. Laine, E. Härkönen, J. Hellsten, J. Lehtinen, and T. Aila. Alias-free generative adversarial networks. *Advances in Neural Information Processing Systems*, 34:852–863, 2021.
- [17] T. Karras, S. Laine, and T. Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019.
- [18] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020.

- [19] J. Ko, K. Cho, D. Choi, K. Ryoo, and S. Kim. 3d gan inversion with pose optimization. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2967–2976, 2023.
- [20] K. Ko, T. Yeom, and M. Lee. Superstargan: Generative adversarial networks for image-to-image translation in large-scale domains. *Neural Networks*, 162:330–339, 2023.
- [21] C. Z. Lin, D. B. Lindell, E. R. Chan, and G. Wetzstein. 3d gan inversion for controllable portrait image animation. *arXiv preprint arXiv:2203.13441*, 2022.
- [22] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, pages 3730–3738, 2015.
- [23] L. Mescheder, A. Geiger, and S. Nowozin. Which training methods for gans do actually converge? In *International conference on machine learning*, pages 3481–3490. PMLR, 2018.
- [24] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021.
- [25] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.
- [26] T. Nguyen-Phuoc, C. Li, L. Theis, C. Richardt, and Y.-L. Yang. Hologan: Unsupervised learning of 3d representations from natural images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7588–7597, 2019.
- [27] M. Niemeyer and A. Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11453–11464, 2021.
- [28] R. Or-El, X. Luo, M. Shan, E. Shechtman, J. J. Park, and I. Kemelmacher-Shlizerman. Stylesdf: High-resolution 3d-consistent image and geometry generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13503–13513, 2022.
- [29] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [30] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- [31] E. Richardson, Y. Alaluf, O. Patashnik, Y. Nitzan, Y. Azar, S. Shapiro, and D. Cohen-Or. Encoding in style: a stylegan encoder for image-to-image translation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2287–2296, 2021.
- [32] D. Roich, R. Mokady, A. H. Bermano, and D. Cohen-Or. Pivotal tuning for latent-based editing of real images. *ACM Transactions on Graphics (TOG)*, 42(1):1–13, 2022.
- [33] S. Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.
- [34] K. Schwarz, Y. Liao, M. Niemeyer, and A. Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *Advances in Neural Information Processing Systems*, 33:20154–20166, 2020.
- [35] S. M. Seitz and C. R. Dyer. Photorealistic scene reconstruction by voxel coloring. *International journal of computer vision*, 35:151–173, 1999.
- [36] Y. Shen, J. Gu, X. Tang, and B. Zhou. Interpreting the latent space of gans for semantic face editing. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9243–9252, 2020.
- [37] V. Sitzmann, J. Thies, F. Heide, M. Nießner, G. Wetzstein, and M. Zollhofer. Deepvoxels: Learning persistent 3d feature embeddings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2437–2446, 2019.
- [38] H. Song, Y. Du, T. Xiang, J. Dong, J. Qin, and S. He. Editing out-of-domain gan inversion via differential activations. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XVII*, pages 1–17. Springer, 2022.
- [39] M. Tancik, V. Casser, X. Yan, S. Pradhan, B. Mildenhall, P. P. Srinivasan, J. T. Barron, and H. Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8248–8258, 2022.

- [40] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- [41] T. Wei, D. Chen, W. Zhou, J. Liao, Z. Tan, L. Yuan, W. Zhang, and N. Yu. Hairclip: Design your hair by text and reference image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18072–18081, 2022.
- [42] J. Xie, H. Ouyang, J. Piao, C. Lei, and Q. Chen. High-fidelity 3d gan inversion by pseudo-multi-view optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 321–331, 2023.
- [43] B. Xu, N. Wang, T. Chen, and M. Li. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853*, 2015.
- [44] Q. Xu, Z. Xu, J. Philip, S. Bi, Z. Shu, K. Sunkavalli, and U. Neumann. Point-nerf: Point-based neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5438–5448, 2022.
- [45] L. Yang, P. Luo, C. Change Loy, and X. Tang. A large-scale car dataset for fine-grained categorization and verification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3973–3981, 2015.
- [46] Y. Yin, K. Ghasedi, H. Wu, J. Yang, X. Tong, and Y. Fu. Nerfinvertor: High fidelity nerf-gan inversion for single-shot real image animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8539–8548, 2023.
- [47] A. Yu, V. Ye, M. Tancik, and A. Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2021.
- [48] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018.
- [49] T. Zhou, R. Tucker, J. Flynn, G. Fyffe, and N. Snavely. Stereo magnification: Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817*, 2018.
- [50] J. Zhu, Y. Shen, D. Zhao, and B. Zhou. In-domain gan inversion for real image editing. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVII 16*, pages 592–608. Springer, 2020.

## A Introduction of Supplementary Material

In this supplemental document, we offer a detailed overview of the various architectural elements within the network – including the feature fields, the neural renderer, and the discriminator, all discussed in Section B. Furthermore, in Section C, we elucidate a quantitative analysis of our ablation study results, underscoring the efficacy of our loss functions during the training phase of the inverter. In conclusion, we bring forth additional qualitative findings on datasets such as CelebA-HQ [15], CompCar [45], and AFHQ [5]. Two novel experimental approaches are also introduced: the style-mixed 3D representation of two facial input images, and the generation of two objects within a single scene using a generator trained on single-object scenes.



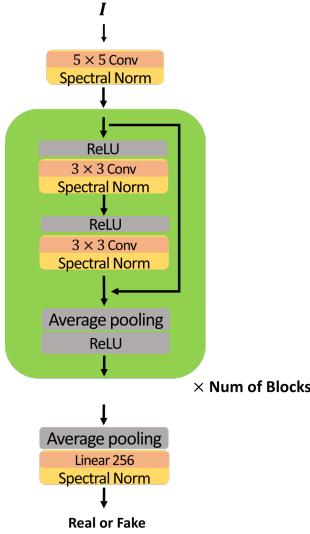
**Figure S8: Architecture of the feature fields and neural renderer.** The feature fields are parameterized with multi-layer perceptrons (MLPs) as shown in the (a). The 3D point  $\mathbf{x}$ , viewing direction  $\mathbf{d}$ , and latent codes  $\mathbf{z}_s$ ,  $\mathbf{z}_a$  are mapped into a volume density  $\sigma$  and feature  $\mathbf{f}$ . In (b), the neural renderer blocks depict the transformation of the volume-rendered feature image  $F$  into final synthesized image  $\hat{I}$ .  $UP_{NN}$  and  $UP_{BL}$  symbolize the nearest neighbour upsampling and bilinear upsampling, respectively.

## B Network Architectures

In this section, we provide the details of network architecture: feature fields, neural renderer, and the discriminator as exhibited in Fig. S8 and S9.

Fig. S8 presents a detailed overview of the architecture underpinning the feature fields and the neural renderer. The construct of the feature fields is parameterized via multi-layer perceptrons, colloquially referred to as MLPs, a feature vividly displayed in subfigure (a). This setup maps a three-dimensional point, the viewing direction, along with latent codes into a volume density and a feature. Subfigure (b) unravels the process behind the neural renderer blocks, demonstrating how these blocks transform a volume-rendered feature image, into the ultimate synthesized image.

Fig. S9 explicates the architecture of the discriminator network, emphasizing the steps involved in processing the input image. Initially, the image is subjected to a series of residual convolution blocks, which are fortified with spectral normalization. This is followed by the execution of an average pooling operation. The process culminates with the derivation of the output probability, which is obtained post the final linear layer, again, involving spectral normalization.



**Figure S9: Architecture of the discriminator.** The input image is processed through residual convolution blocks fortified with spectral normalization, and an average pooling operation. The output probability is derived after the final linear layer with spectral normalization.

| Ablation Losses   | FID          |
|-------------------|--------------|
| $L_{latent}$      | 80.08        |
| $+L_{reconst}$    | 17.82        |
| $+L_{GAN}$        | 15.53        |
| <b>Full model</b> | <b>14.77</b> |

**Table S2: FID score of the ablation study.** The full model has been trained with latent loss, reconstruction loss, GAN loss, and perceptual loss.

## C Supplementary Experimental Results

### C.1 The Necessity of Loss Components in Training Session: A Quantitative Evaluation

Tab. S2 offers a quantitative testament to the indispensable nature of the loss components used in the training session of the inverter. These encompass latent loss, reconstruction loss, GAN loss, and perceptual loss. It is observed that the Fréchet Inception Distance (FID) [12] experiences a steady enhancement with each loss component incrementally added to the naive model, which originally only employs the latent loss.

### C.2 Extended Operational Results

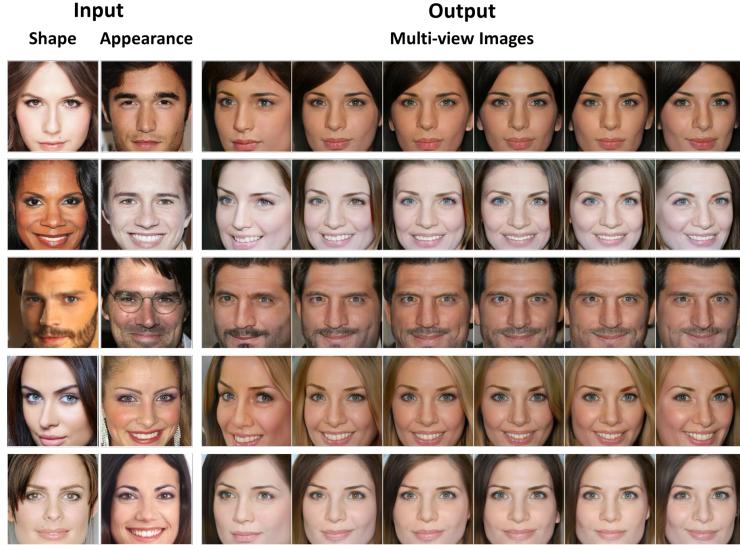
In this section, we present the application results of the proposed model through Fig. S10 and S11, showcasing style-mixed 3D synthesis and the generation of two objects within a single scene.

Our model demonstrates a unique ability to generate multiple objects within a single scene, even when trained on a dataset consisting primarily of single-object scenes. This is accomplished by leveraging multiple decoder segments within our network architecture. Although our empirical exploration has only been executed on one dataset, the theoretical underpinnings suggest a promising generalizability of this phenomenon. A testament to the robustness of our model is its successful exhibition of zero-shot learning capabilities, as evidenced by an experiment where two CompCars images are synthesized into one image. Like the generation of individual objects, each object within the composite scene retains the ability to undergo transformations such as longitudinal displacement and rotation.

Additionally, we incorporate style mixing in our model with the application of the inverter structure we proposed, utilizing the CelebA-HQ dataset. In the style mixing paradigm that we suggest, our inverter, producing two distinct outputs, generates a shape vector from one image, and an appearance vector from another. These two vectors are subsequently utilized as input for the generator to synthesize a novel object. This process further underscores the model’s zero-shot learning capability.



**Figure S10: Generating two objects in a single scene.** Results exhibit the compositional scene representation by generating two objects in a single scene. The inverter transforms two input images into two sets of the latent codes, and the generator which trained on single-object scenes synthesizes a single scene including two independent objects.



**Figure S11: Multi-view images with style mixing of two CelebA-HQ input images.** The invertor extracts the latent codes from two independent input images for generating style mixed images. Each output object is generated by  $\mathbf{z}_s$  of the first image and  $\mathbf{z}_a$  of the second image.

### C.3 Supplementary Results

Fig. S12, S13, and S14 deliver additional examples on CelebA-HQ [15], AFHQ [5], and CompCar [45] datasets.

We embark on rigorous evaluation of our model using a diverse range of input images sourced from varied datasets. With the CelebA dataset, we assess the model’s performance using faces of different genders, ages, and ethnic backgrounds, all of which yield impressive quality in output. In the context of the AFHQ dataset, we utilize images from a variety of categories as input for our testing phase. It

is worth noting that these results, encompassing distinct categories, are obtained using a single model with different conditional vector inputs, thereby highlighting the large capacity of our model.

The CompCars dataset allows us to experiment with 360-degree image generation using real image inputs representing various car models, colours, and camera poses. It is important to note that a significant advantage of our model is the freedom it provides in the longitudinal movement of objects, along with the capacity to alter the background. This flexibility underpins the model's capacity for highly controllable image synthesis, an attribute that holds immense potential for a wide array of applications.

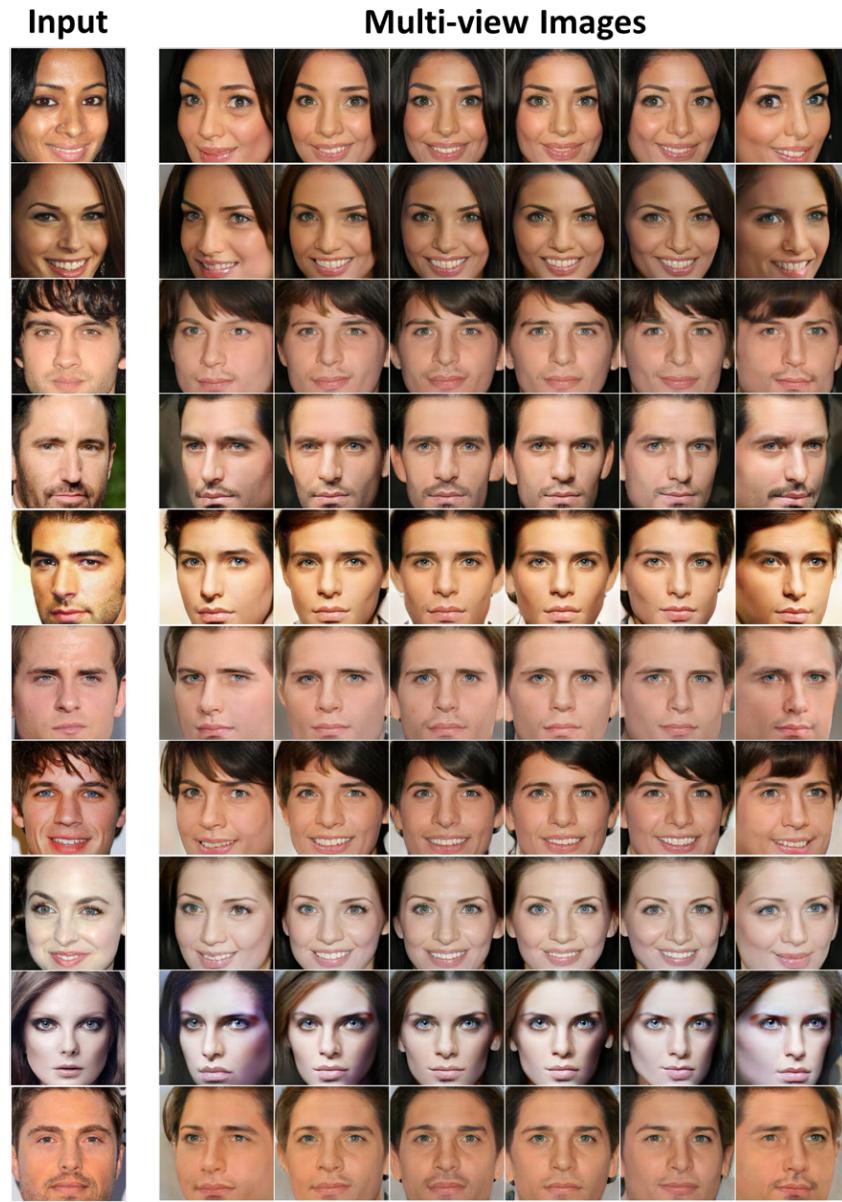


Figure S12: **Supplementary results with  $256^2$  CelebA-HQ image inputs.**

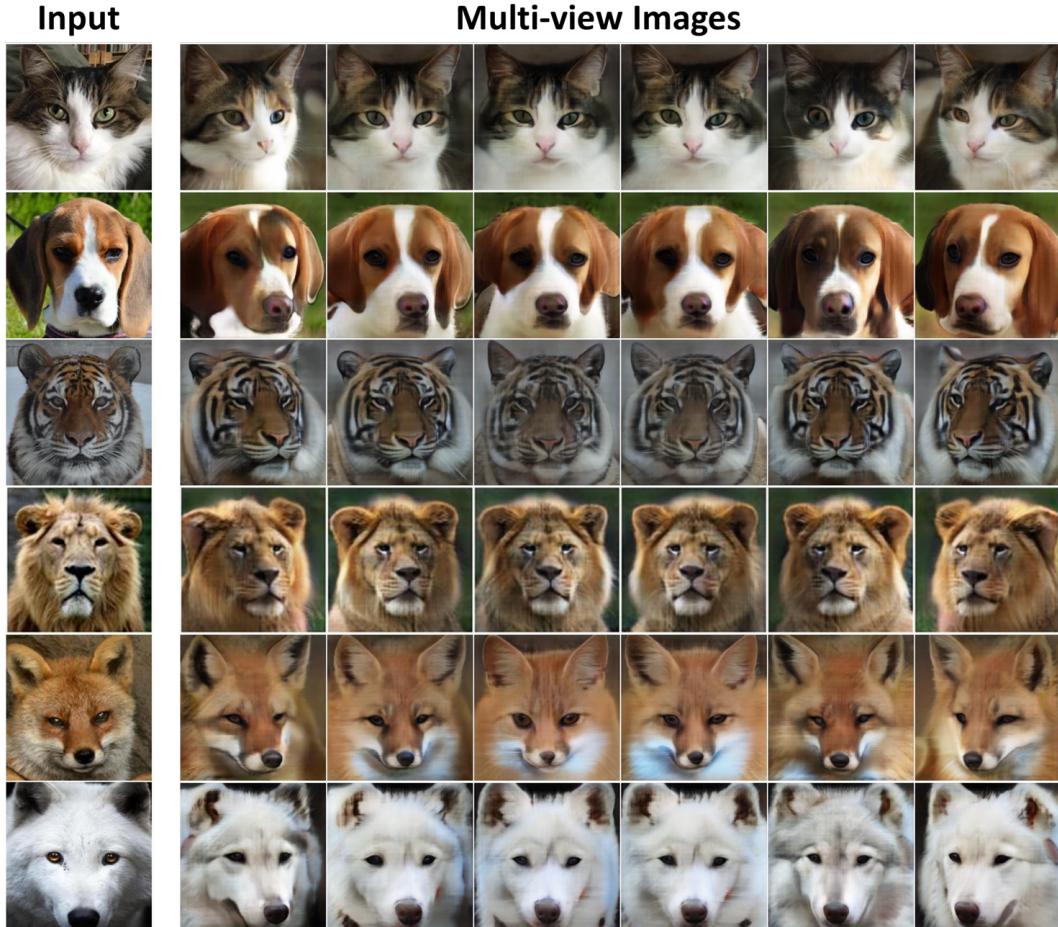


Figure S13: Supplementary results with  $256^2$  AFHQ image inputs.

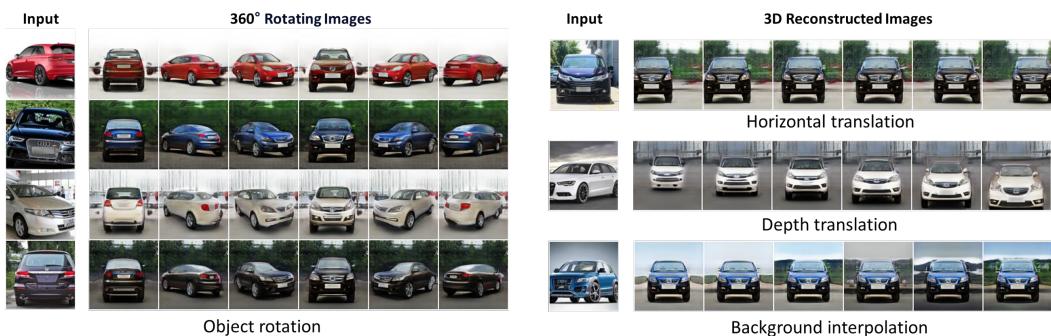


Figure S14: Controllable image synthesis with  $128^2$  CompCars image inputs.