

CLIP-NeRF: Text-and-Image Driven Manipulation of Neural Radiance Fields

2022.03.14

Presenter: Junha Hyung

Introduction

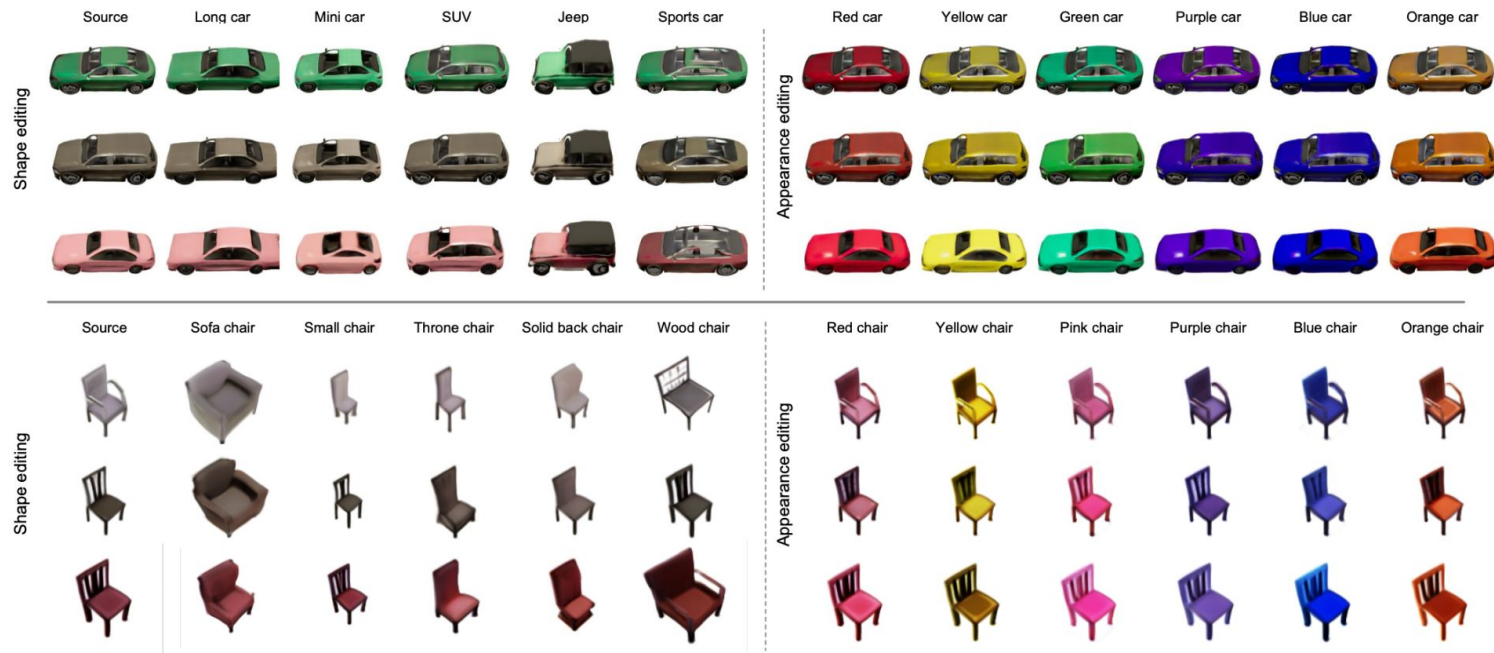


Figure 6. Text-Driven Editing Results.

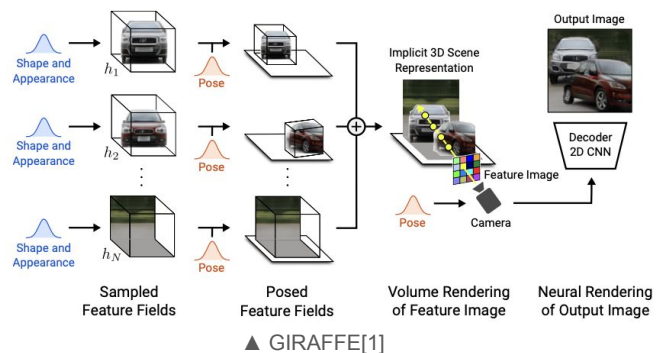
Contribution

- The first text-and-image-driven manipulation method for NeRF, using a unified framework to provide users with flexible control over 3D content using either a text prompt or an exemplar image.
- A disentangled conditional NeRF architecture by introducing a shape code to deform the volumetric field and an appearance code to control the emitted colors.
- Feedforward code mappers that enable the fast inference

Related works

Generative Neural Radiance Field

- object와 background의 shape/appearance code를 sampling해서, arbitrary한 pose에 따라 image들을 rendering.
- 기존 2D-GANs와는 달리, camera pose를 조작할 수 있음.



Edit NeRF

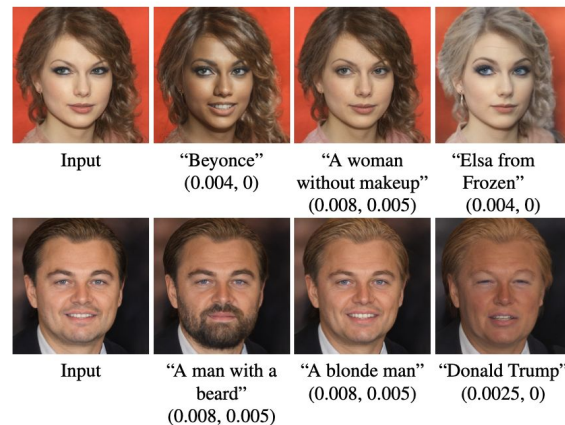
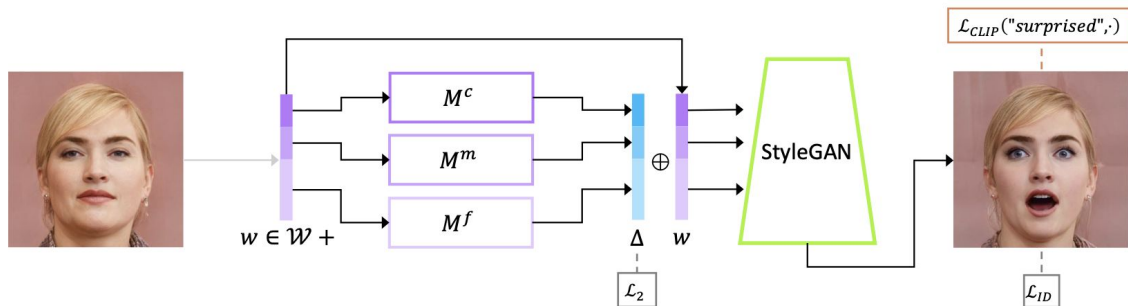
- Scribble로 color, shape editing
- Optimization based



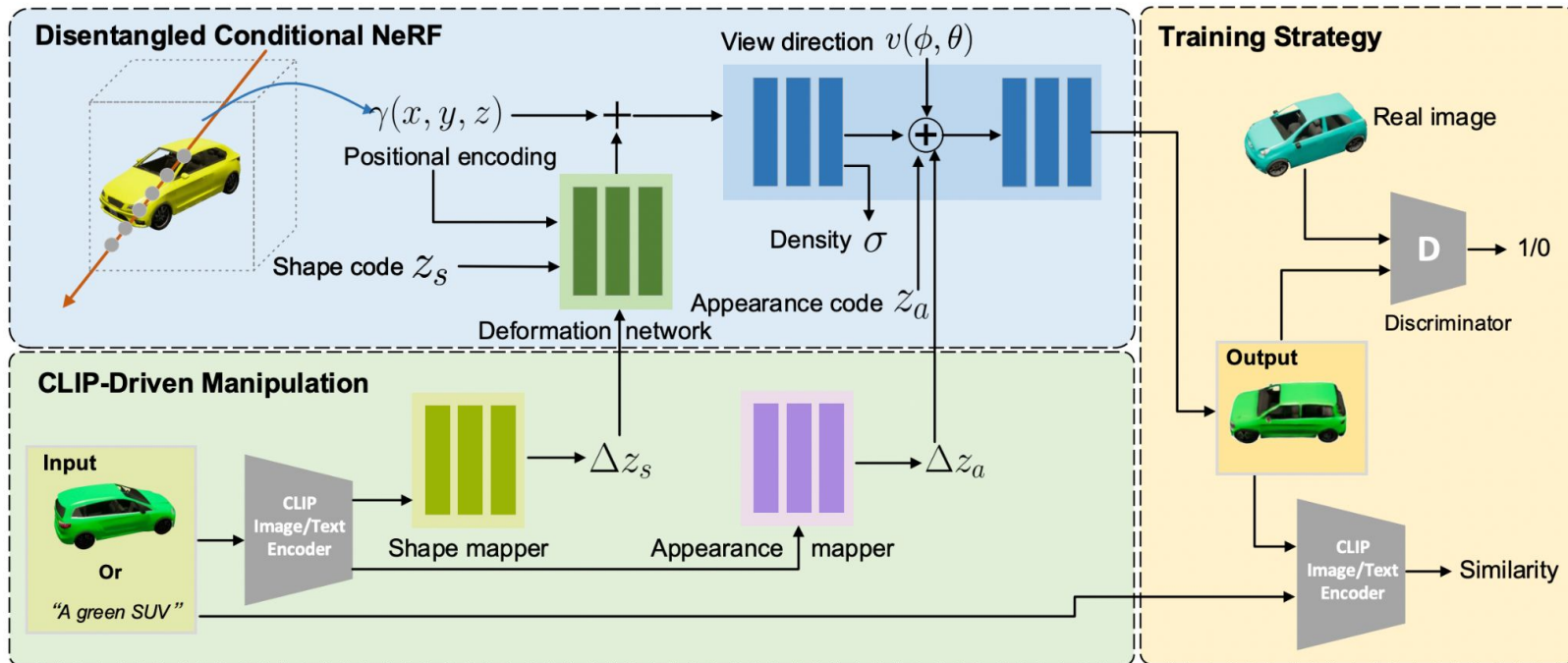
Related works

Style Clip [1]

- Pretrained StyleGAN과, Clip 활용해 text guided image editing



Method



Overall formulation

$$\mathcal{F}'_{\theta}(\mathbf{x}, \mathbf{v}, \mathbf{z}_s, \mathbf{z}_a) : (\Gamma(\mathbf{x}) \oplus \mathbf{z}_s, \Gamma(\mathbf{v}) \oplus \mathbf{z}_a) \rightarrow (\mathbf{c}, \sigma), \quad (1)$$

Positional encoding

$$\gamma(p)_k = \begin{cases} \sin(2^k \pi p), & \text{if } k \text{ is even,} \\ \cos(2^k \pi p), & \text{if } k \text{ is odd,} \end{cases}$$

Conditional shape deformation

- Completely isolates the shape condition from affecting the appearance
- Regularize the output shape to be smooth deformation of the base shape

$$\gamma^*(p, \Delta p)_k = \gamma(p)_k + \tanh(\Delta p_k),$$

$$p \in \mathbf{p}, \Delta p \in \mathcal{T}(\mathbf{p}, \mathbf{z}_s)$$

Clip mapper

$$\mathbf{z}_s = \mathcal{M}_s(\hat{\mathcal{E}}_t(\mathbf{t})) + \mathbf{z}'_s,$$

$$\mathbf{z}_a = \mathcal{M}_a(\hat{\mathcal{E}}_t(\mathbf{t})) + \mathbf{z}'_a,$$

Training Strategy

- Two phase training
 - First train generative NeRF
 - Second freeze NeRF and train mapper networks

$$\begin{aligned} \mathcal{L}_{\text{GAN}} = & \mathbb{E}_{\mathbf{z}_s \sim \mathcal{Z}_s, \mathbf{z}_a \sim \mathcal{Z}_a, \mathbf{v} \sim \mathcal{Z}_v} [f(\mathcal{D}(\mathcal{F}_\theta(\mathbf{v}, \mathbf{z}_s, \mathbf{z}_a)))] \\ & + \mathbb{E}_{\mathbf{I} \sim d} [f(-\mathcal{D}(\mathbf{I}) + \lambda_r \|\nabla \mathcal{D}(\mathbf{I})\|^2)]. \end{aligned} \quad (7)$$

$$D_{\text{CLIP}}(\mathbf{I}, \mathbf{t}) = 1 - \langle \hat{\mathcal{E}}_i(\mathbf{I}), \hat{\mathcal{E}}_t(\mathbf{t}) \rangle,$$

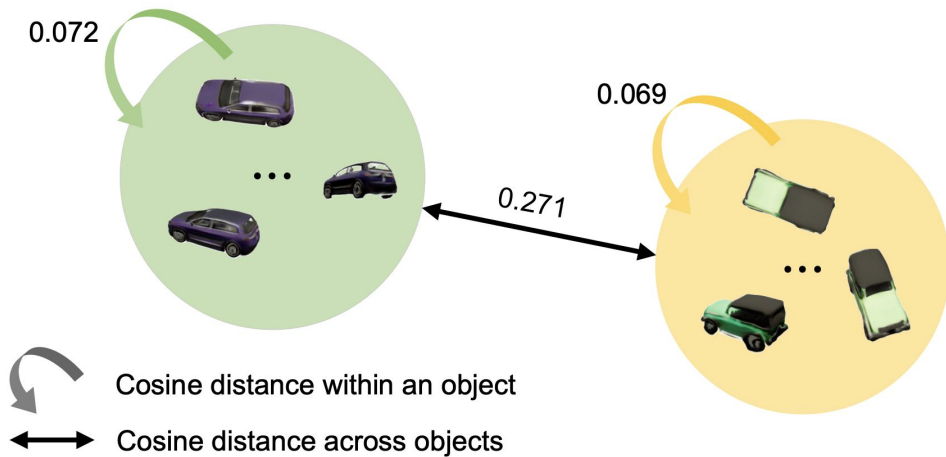
$$\begin{aligned} \mathcal{L}_{\text{shape}} = & f(\hat{\mathcal{D}}(\hat{\mathcal{F}}_\theta(\mathbf{v}, \mathcal{M}_s(\hat{\mathcal{E}}_t(\mathbf{t})) + \mathbf{z}_s, \mathbf{z}_a))) + \\ & \lambda_c D_{\text{CLIP}}(\hat{\mathcal{F}}_\theta(\mathbf{v}, \mathcal{M}_s(\hat{\mathcal{E}}_t(\mathbf{t})) + \mathbf{z}_s, \mathbf{z}_a), \mathbf{t}), \end{aligned} \quad (8)$$

$$\begin{aligned} \mathcal{L}_{\text{appear}} = & f(\hat{\mathcal{D}}(\hat{\mathcal{F}}_\theta(\mathbf{v}, \mathbf{z}_s, \mathcal{M}_a(\hat{\mathcal{E}}_t(\mathbf{t})) + \mathbf{z}_a))) + \\ & \lambda_c D_{\text{CLIP}}(\hat{\mathcal{F}}_\theta(\mathbf{v}, \mathbf{z}_s, \mathcal{M}_a(\hat{\mathcal{E}}_t(\mathbf{t})) + \mathbf{z}_a), \mathbf{t}). \end{aligned} \quad (9)$$

Discussion

Image editing을 하기에 Clip이 적절한 모델인가?

- Higher cosine similarity within an object



Inverse Manipulation

To be specific, during each iteration, we first optimize \mathbf{v} while keeping \mathbf{z}_s and \mathbf{z}_a fixed using the following loss:

$$\mathcal{L}_v = \|\hat{\mathcal{F}}_\theta(\mathbf{v}, \hat{\mathbf{z}}_s, \hat{\mathbf{z}}_a) - \mathbf{I}_r\|_2 + \lambda_v D_{\text{CLIP}}(\hat{\mathcal{F}}_\theta(\mathbf{v}, \hat{\mathbf{z}}_s, \hat{\mathbf{z}}_a), \mathbf{I}_r). \quad (10)$$

We then update the shape code by minimizing:

$$\mathcal{L}_s = \|\hat{\mathcal{F}}_\theta(\hat{\mathbf{v}}, \mathbf{z}_s + \lambda_n \mathbf{z}_n, \hat{\mathbf{z}}_a) - \mathbf{I}_r\|_2 + \lambda_s D_{\text{CLIP}}(\hat{\mathcal{F}}_\theta(\hat{\mathbf{v}}, \mathbf{z}_s + \lambda_n \mathbf{z}_n, \hat{\mathbf{z}}_a), \mathbf{I}_r),$$

$$\mathcal{L}_a = \|\hat{\mathcal{F}}_\theta(\hat{\mathbf{v}}, \hat{\mathbf{z}}_s, \mathbf{z}_a + \lambda_n \mathbf{z}_n) - \mathbf{I}_r\|_2 + \lambda_a D_{\text{CLIP}}(\hat{\mathcal{F}}_\theta(\hat{\mathbf{v}}, \hat{\mathbf{z}}_s, \mathbf{z}_a + \lambda_n \mathbf{z}_n), \mathbf{I}_r),$$

$\hat{\mathbf{v}}$ fixed, \mathbf{z}_n is a random standard Gaussian added in each step to improve the optimization and λ_n linearly decays from 1 to 0 through

Experiments

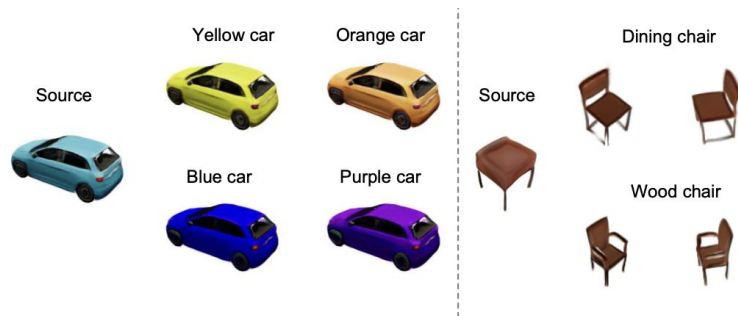
Dataset

- CARLA(10k cars, 256x256)
- Photoshapes(150k chairs, 128x128)

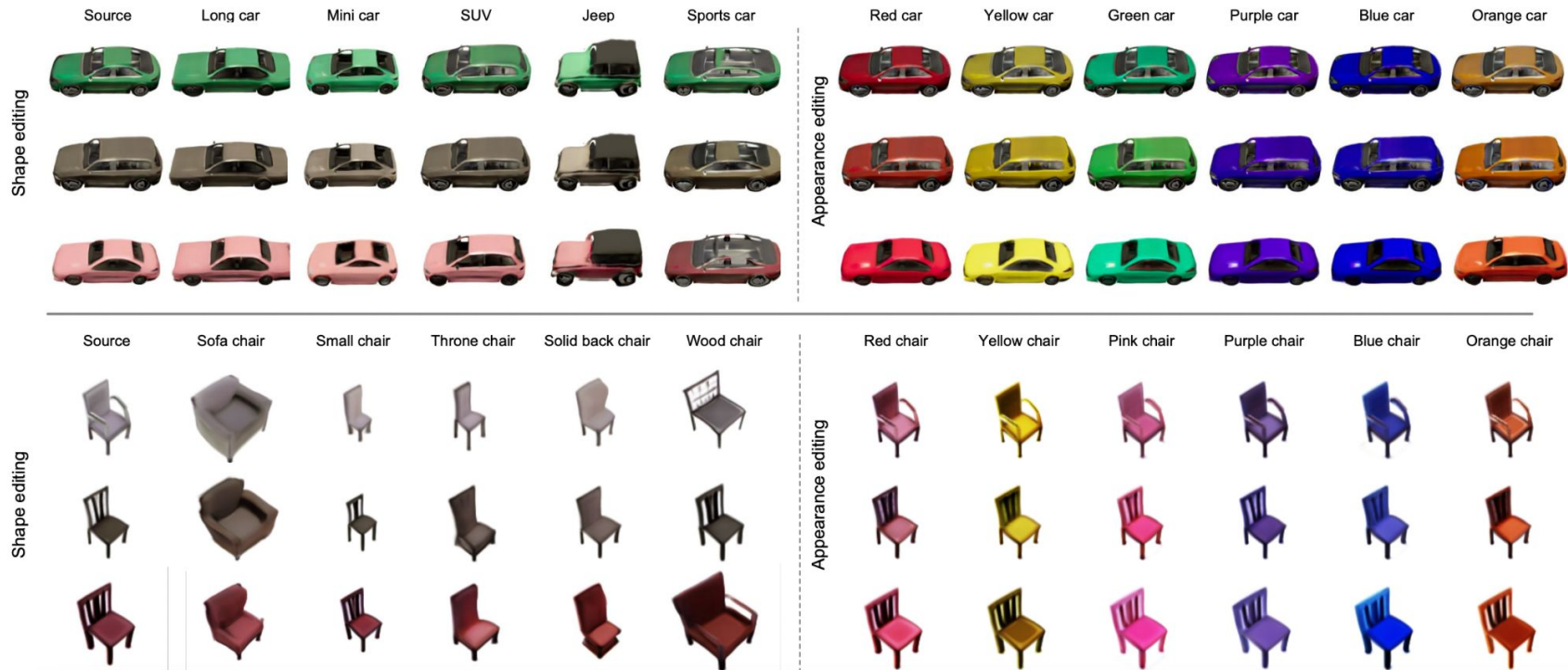
Inference time

	Chairs		Cars	
	Shape	Appearance	Shape	appearance
EditNeRF	30.0	15.9	33.2	16.8
Ours	0.58	0.51	2.12	1.98

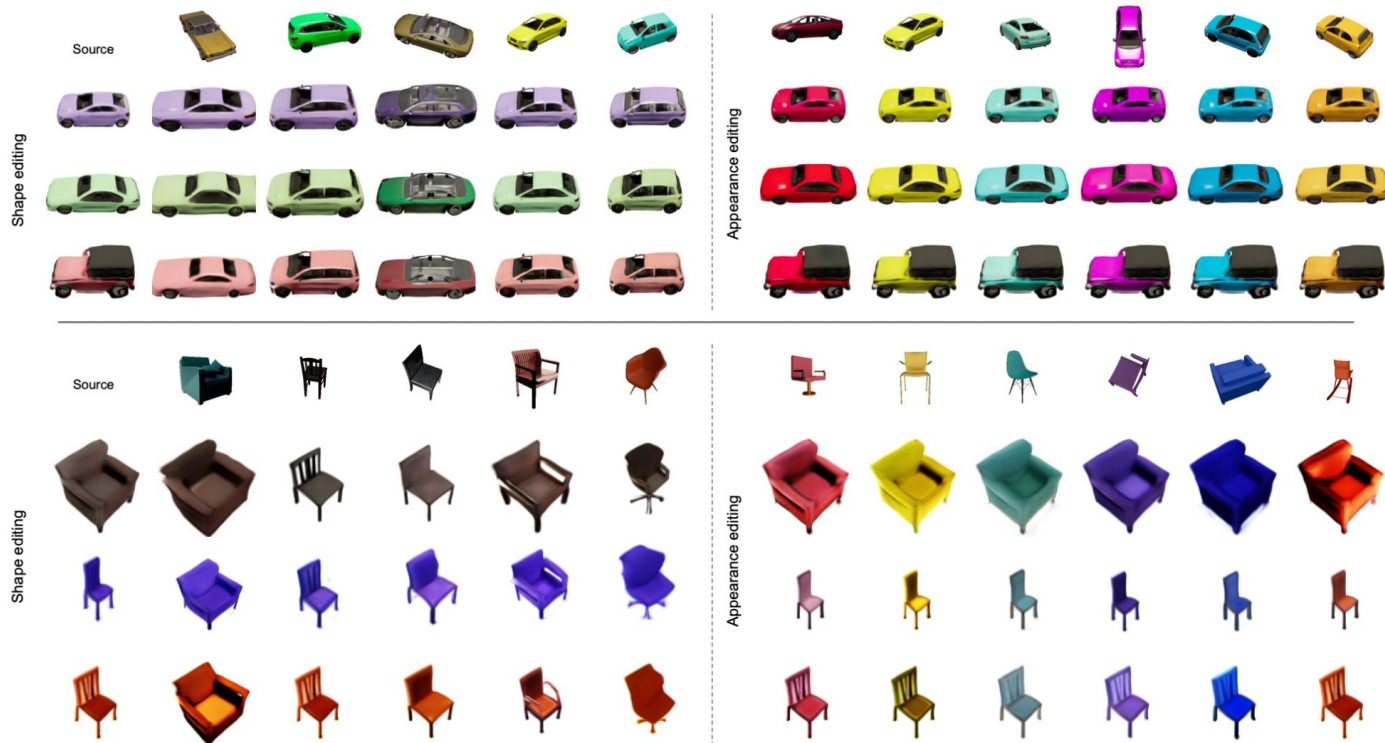
Table 1. **Compared to EditNeRF [21] on editing time averaged on 20 images.** We only include the inference/optimization time(s) and single-view rendered time(s) for chairs (128×128 pixels) and cars (256×256 pixels).



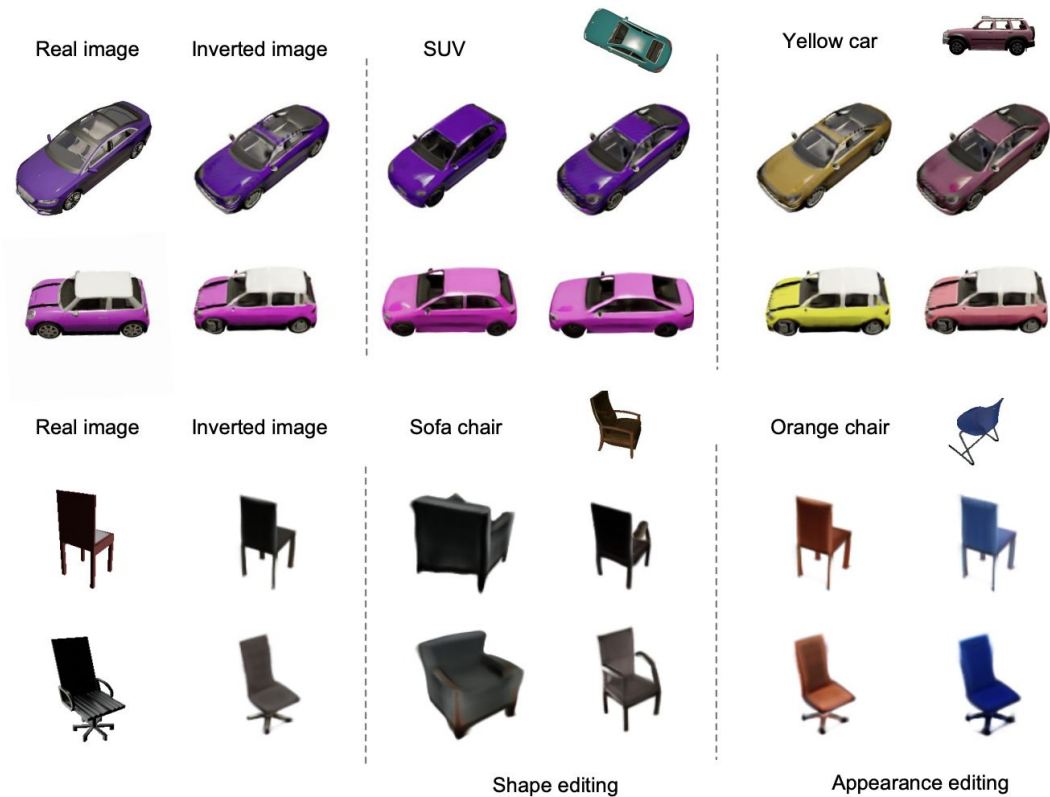
Results - text driven



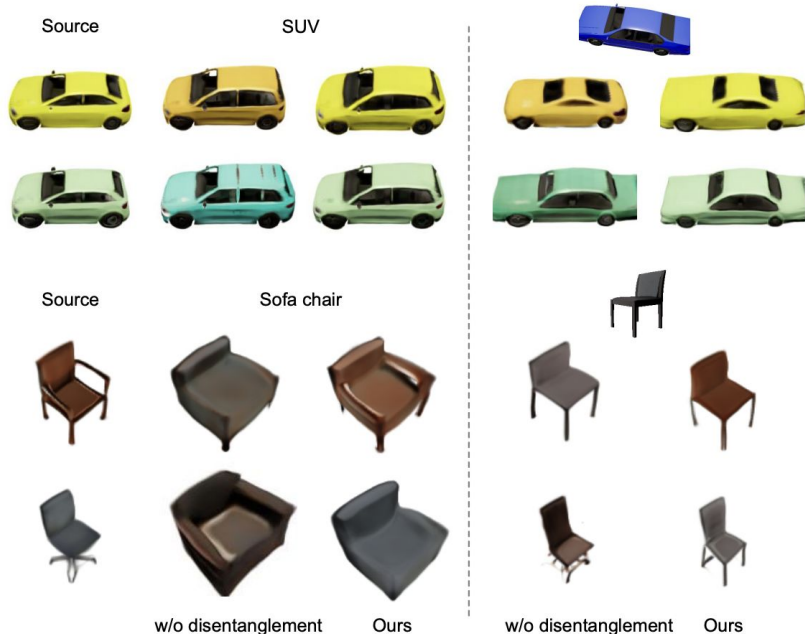
Results - exemplar driven



Results - Editing Real Images



Ablation



	Chairs			Cars		
	Before	After	Diff.	Before	After	Diff.
EditNeRF	36.8	40.2	3.4	102.8	118.7	15.9
(a) w/o disen.	52.5	54.3	1.8	69.2	69.9	0.7
Ours	47.8	49.0	1.2	66.7	67.2	0.5
(b) w/o disen.	52.5	53.2	0.7	69.2	71.1	1.9
Ours	47.8	48.4	0.6	66.7	67.8	1.1

Table 2. **Fréchet inception distance (FID) for evaluating the image quality of reconstructed views before and after editing on: (a) color and (b) shape (lower value means better).** We use 2K images with various views drawn randomly from the latent space to calculate the FIDs for reconstructed images, and then perform various edits on these images to recalculate FIDs of edited results.

Ablation

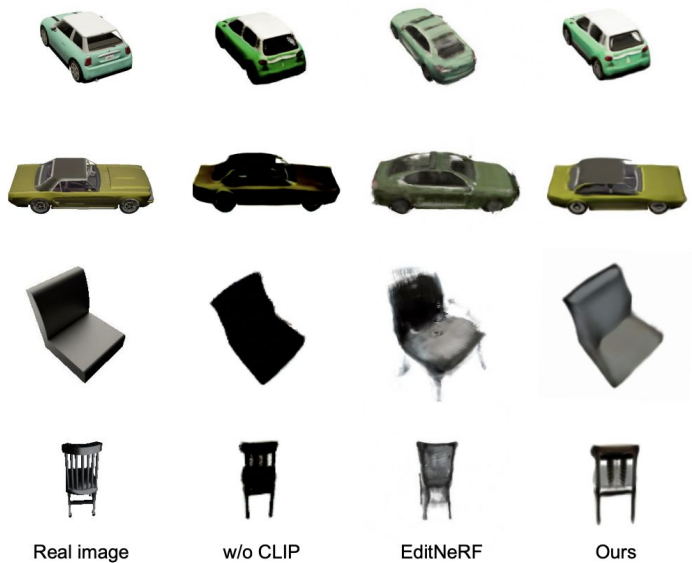


Figure 5. Ablation study on our inversion method and comparison with EditNeRF.

Limitations

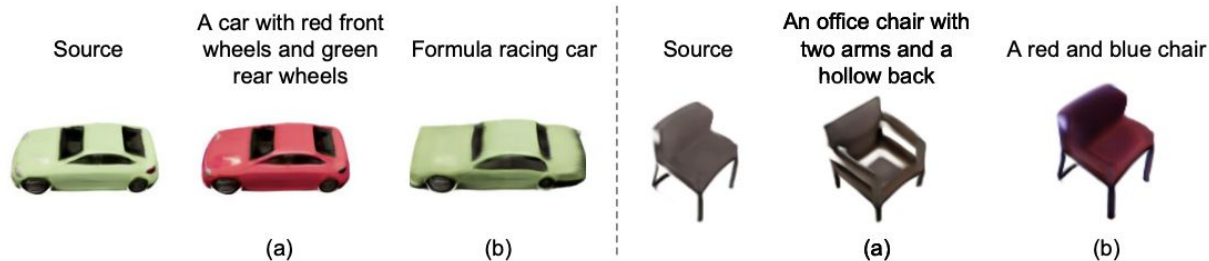


Figure 11. **Limitations.** Our method cannot handle fine-grained edits (a) and out-of-domain edits (b).