

CLIP-NeRF:

Text-and-Image Driven Manipulation of Neural Radiance Fields

Can Wang
City University of Hong Kong
cwang355-c@my.cityu.edu.hk

Menglei Chai
Snap Inc.
cmlatsim@gmail.com

Mingming He
USC Institute for Creative Technologies
hmm.lillian@gmail.com

Dongdong Chen
Microsoft Cloud AI
cddlyf@gmail.com

Jing Liao*
City University of Hong Kong
jingliao@cityu.edu.hk

target

- 对于NeRF的condition-driven manipulation,主要是改变shape和appearance

问题

- 形状和颜色相互联系，另外一个会随着其中一个改变
- 想要解耦appearance 和 shape

Method

$$F : (\mathbf{x}, \mathbf{v}) \rightarrow (\mathbf{c}, \sigma)$$

$$F_\theta : (\mathbf{x}, \mathbf{v}, \mathbf{z_s}, \mathbf{z_a}) \rightarrow (\mathbf{c}, \sigma)$$

$$(\Gamma(\mathbf{x}) \oplus z_s, \Gamma(\mathbf{v}) \oplus z_a) \rightarrow (\mathbf{c}, \sigma)$$

Solution

Shape Deformation Network τ

$$\begin{aligned}\tau : (\mathbf{x}, z_s) &\rightarrow \Delta\mathbf{x} \\ \Gamma^*(\mathbf{p}, z_s) &= \{\gamma^*(p, \Delta p) | p \in \mathbf{p}, \Delta p \in \tau(\mathbf{p}, z_s)\} \\ \gamma^*(p, \Delta p)_k &= \gamma(p)_k + \tanh(\Delta p_k)\end{aligned}$$

Deferred Appearance Conditioning

$$\begin{aligned}\mathcal{F}_\theta(\mathbf{x}, \mathbf{v}, z_s, z_a) : (\Gamma^*(\mathbf{x}, z_s), \Gamma(\mathbf{v}) \oplus z_a) &\rightarrow (\mathbf{x}, \sigma) \\ \mathcal{F}_\theta(\mathbf{v}, z_s, z_a), x \in R\end{aligned}$$

CLIP-Driven Manipulation

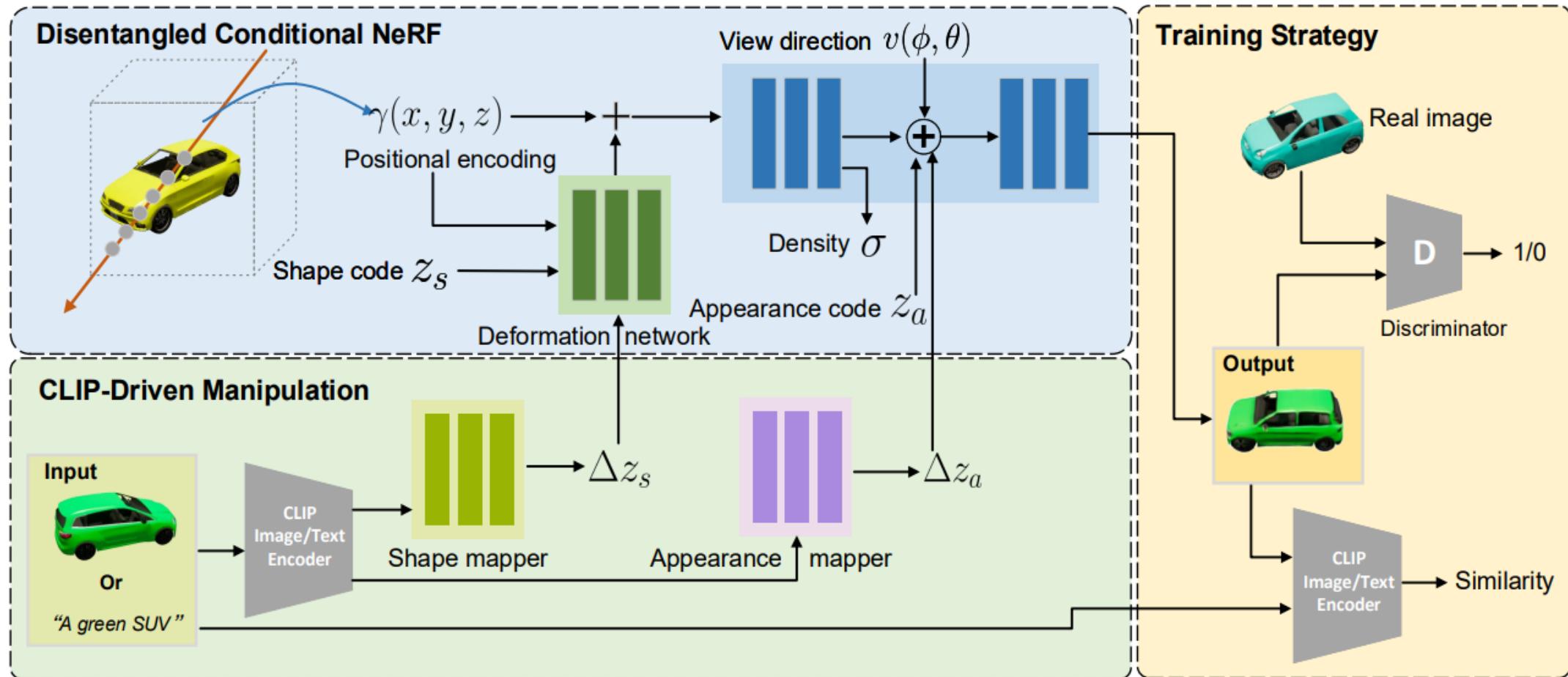
initial state: z'_s, z'_a

trained Mapper: $\mathcal{M}_s, \mathcal{M}_a$

$$\begin{aligned} z_s &= \mathcal{M}_s(\hat{\xi}_t(\mathbf{t})) + z'_s, \\ z_a &= \mathcal{M}_a(\hat{\xi}_t(\mathbf{t})) + z'_a, \end{aligned}$$

cross-modal CLIP distance function:

$$D_{CLIP}(I, t) = 1 - \langle \hat{\xi}_i(I), \hat{\xi}_t(t) \rangle$$



training strategy

two stages:

- disentangled conditional NeRF including the conditional NeRF generator(\mathcal{F}) and the deformation network(τ);
- fix the weights of the generator(\mathcal{F}) and train the CLIP manipulation parts including both the shape and appearance mappers. ($\mathcal{M}_s, \mathcal{M}_a$)

Disentangled Conditional NeRF. Our conditional NeRF generator \mathcal{F}_θ is trained together with the deformation network using a non-saturating GAN objective [22] with the discriminator \mathcal{D} , where $f(x) = -\log(1 + \exp(-x))$ and λ_r is the regularization weight. Assuming that real images \mathbf{I} form the training data distribution of d , we randomly sample the shape code \mathbf{z}_s , the appearance code \mathbf{z}_a , and the camera pose from \mathcal{Z}_s , \mathcal{Z}_a , and \mathcal{Z}_v , respectively, where \mathcal{Z}_s and \mathcal{Z}_a are the normal distribution, and \mathcal{Z}_v is the upper hemisphere of the camera coordinate system.

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

CLIP Manipulation Mappers. We use pre-trained NeRF generator \mathcal{F}_θ , CLIP encoders $\{\hat{\mathcal{E}}_t, \hat{\mathcal{E}}_i\}$, and the discriminator \mathcal{D} to train the CLIP shape mapper \mathcal{M}_s and appearance mapper \mathcal{M}_a . All network weights, except the mappers, are fixed, denoted as $\{\cdot\}$. Similar to the first stage, we randomly sample the shape code \mathbf{z}_s , the appearance code \mathbf{z}_a , and the camera pose \mathbf{v} from their respective distributions. In addition, we sample the text prompt \mathbf{t} from a pre-defined text library \mathbf{T} . By using our CLIP distance D_{CLIP} (Eq. 6) with weight λ_c , we train the mappers with the following losses:

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

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

Inverse Manipulation

获得初始参数，采用迭代优化法

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:

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

We then update the shape code by minimizing:

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

The appearance code is updated in a similar manner:

$$\begin{aligned}\mathcal{L}_a = & \|\hat{\mathcal{F}}_\theta(\hat{\mathbf{v}}, \hat{\mathbf{z}}_s, \mathbf{z}_a + \lambda_n \mathbf{z}_n) - \mathbf{I}_r\|_2 + \\ & \lambda_a D_{\text{CLIP}}(\hat{\mathcal{F}}_\theta(\hat{\mathbf{v}}, \hat{\mathbf{z}}_s, \mathbf{z}_a + \lambda_n \mathbf{z}_n), \mathbf{I}_r),\end{aligned} \quad (12)$$

Experiments

datasets

- Photoshapes (chairs)
- Carla (cars)

comparision

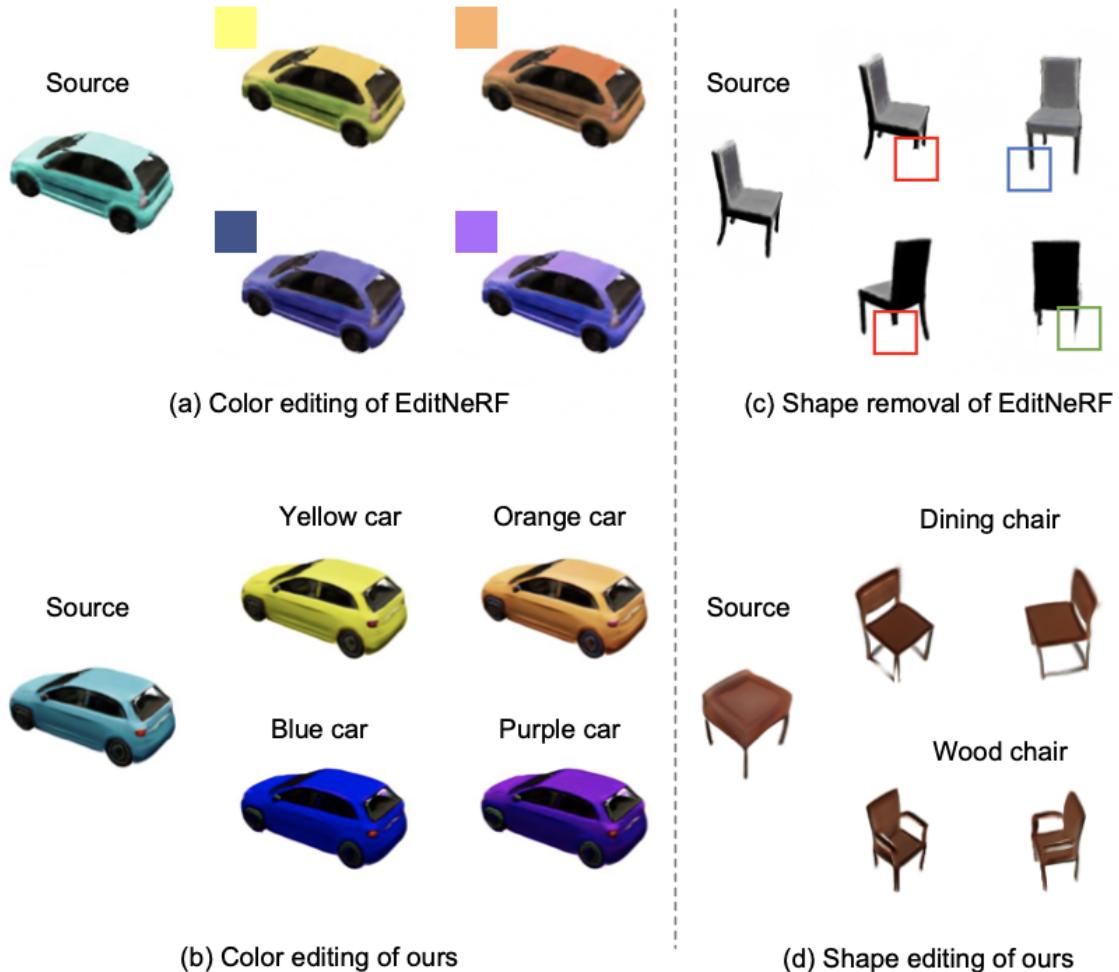


Figure 3. Compared to EditNeRF.

	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).

Ablation

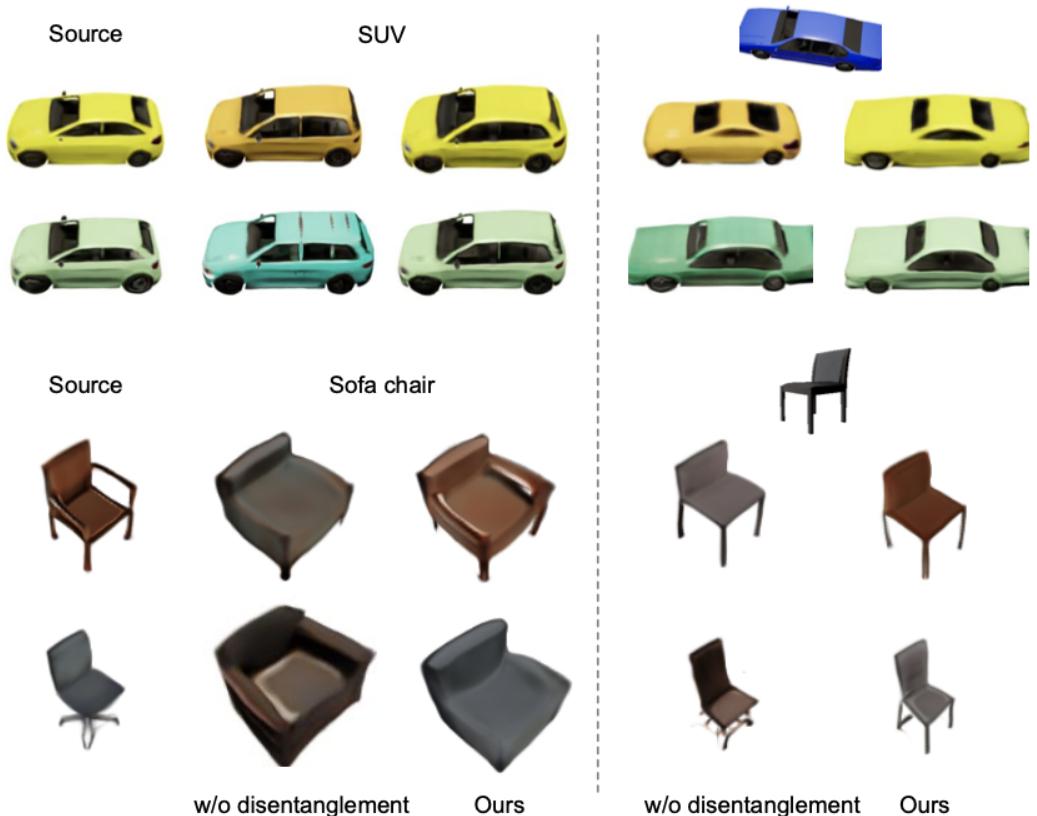


Figure 4. Ablation Study for Disentanglement. We show text-and-exemplar driven shape editing results of our method and the baseline method without using our disentangled technique. When editing the shape, the latter can change the appearance, while ours keeps the appearance unchanged.

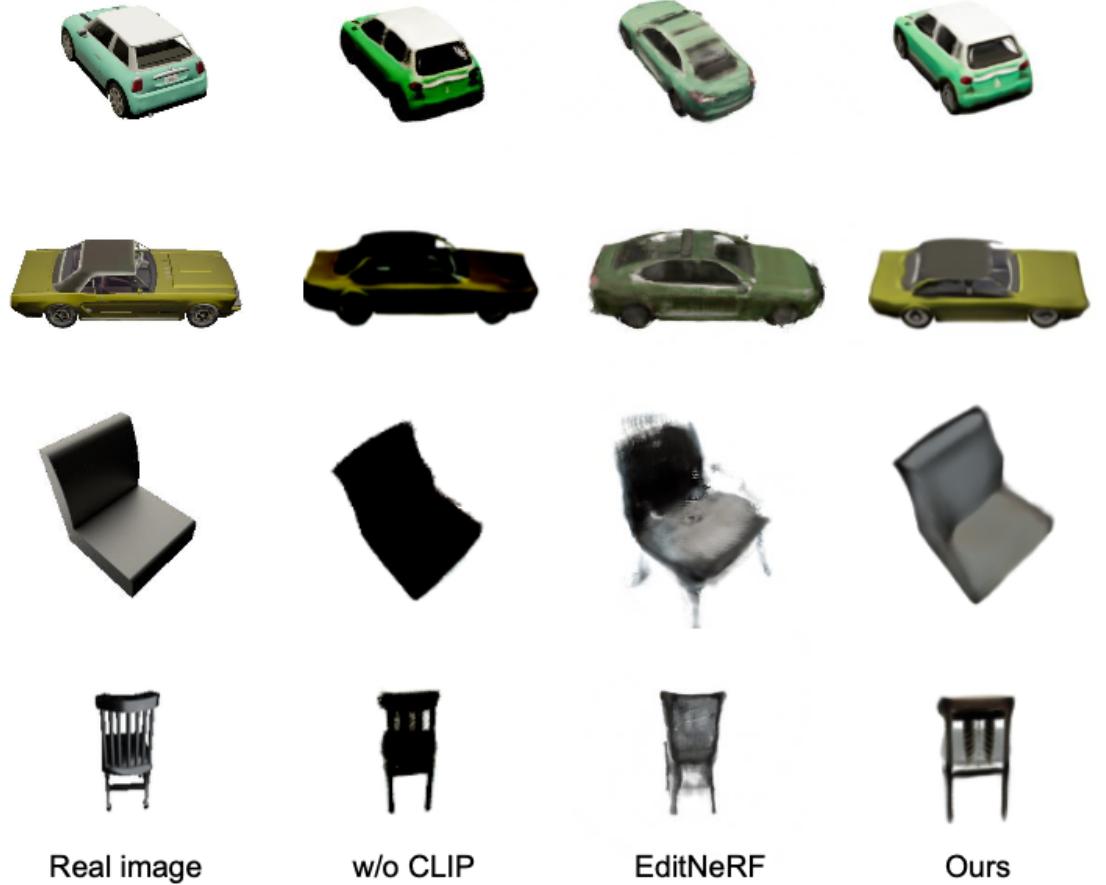


Figure 5. Ablation study on our inversion method and comparison with EditNeRF.

	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

Results

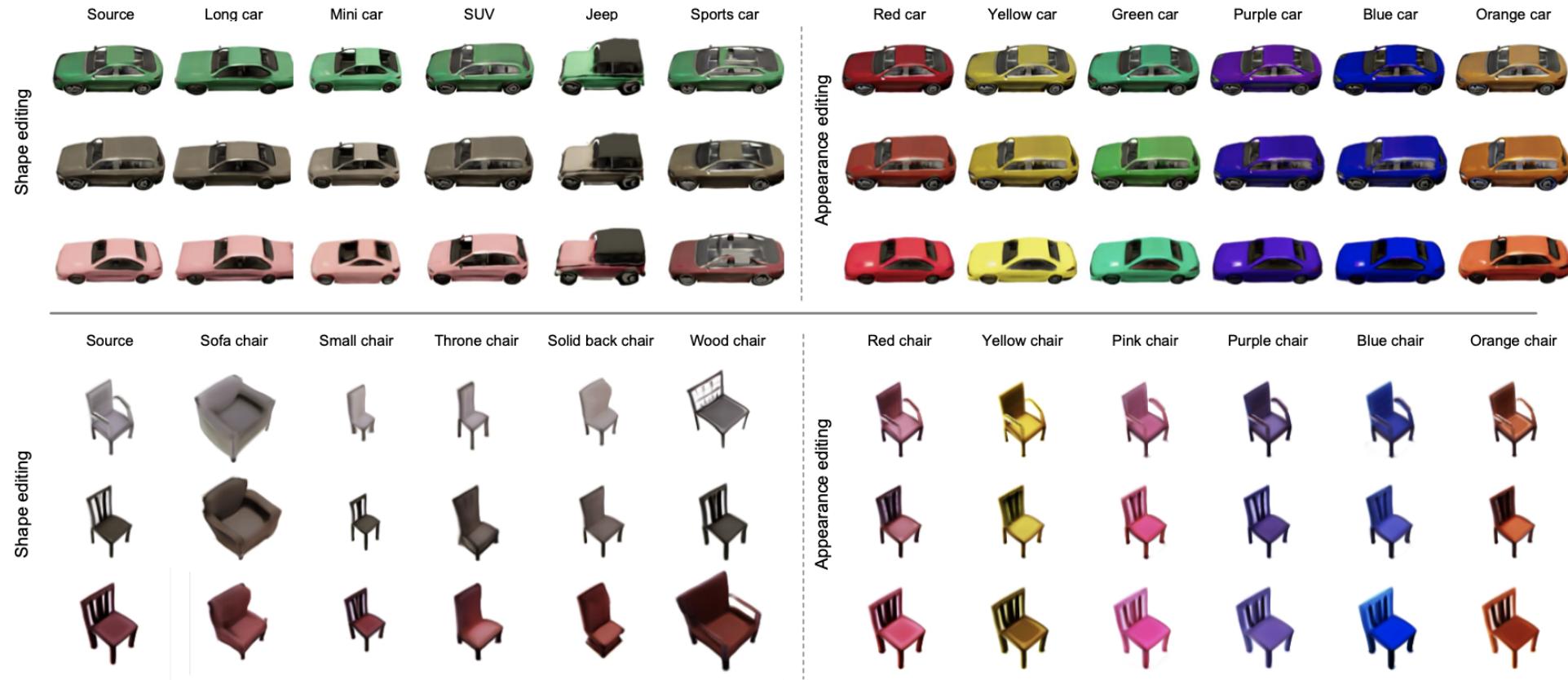


Figure 6. Text-Driven Editing Results.

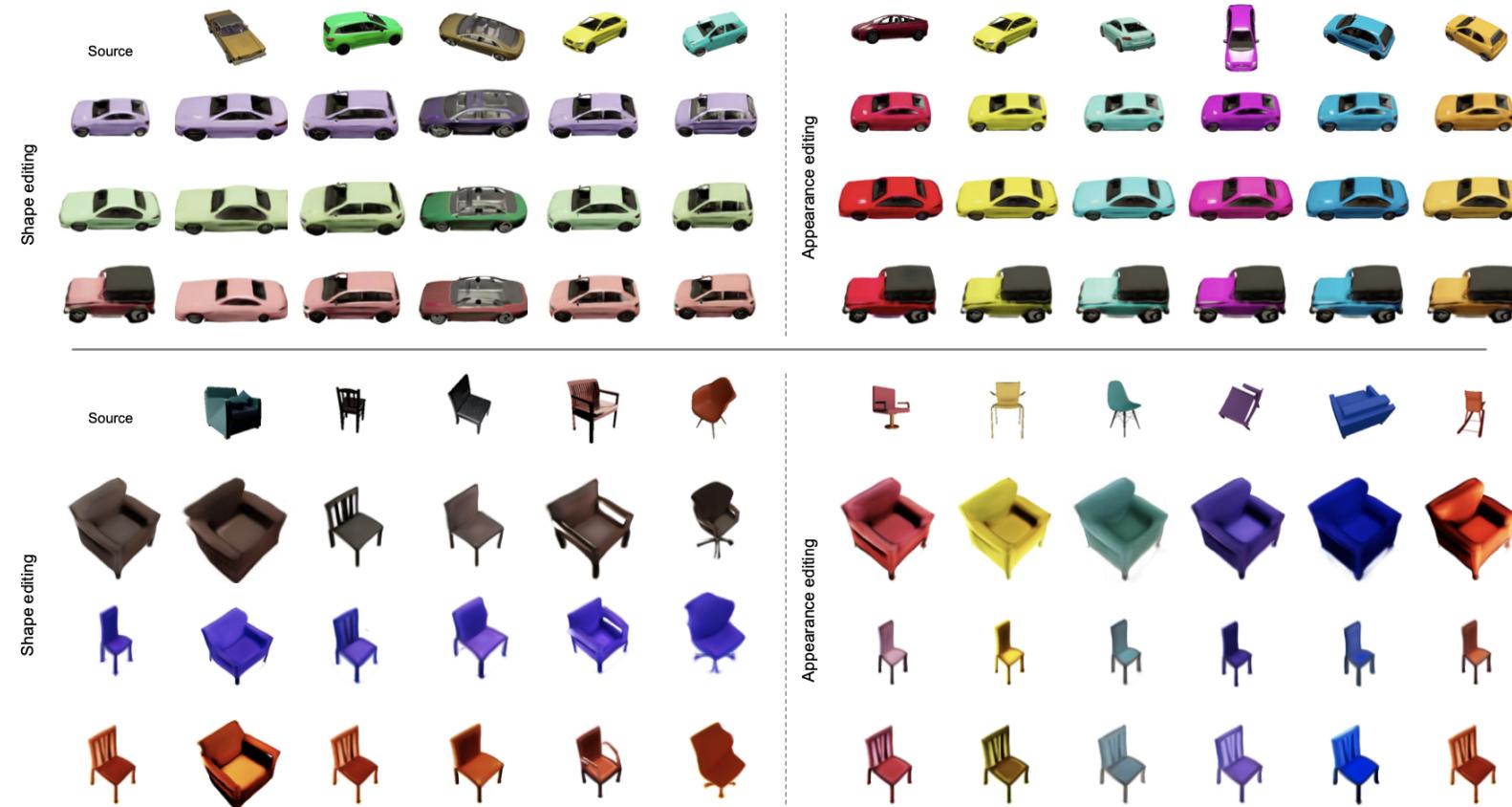


Figure 7. Exemplar-Driven Editing Results.

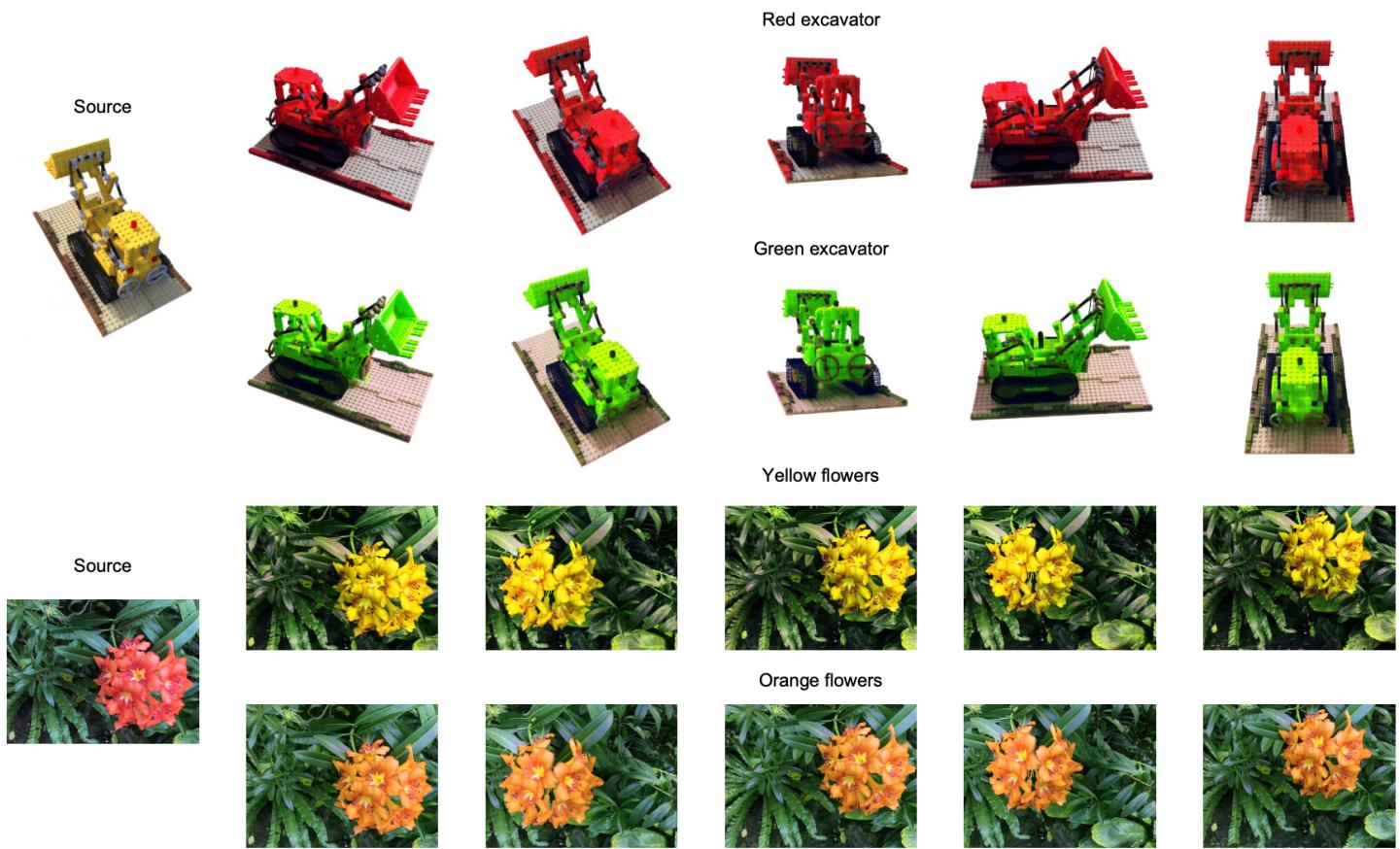


Figure 14. Single NeRF appearance editing results with our designed CLIP loss. NeRF models are trained on LLFF dataset [23].