



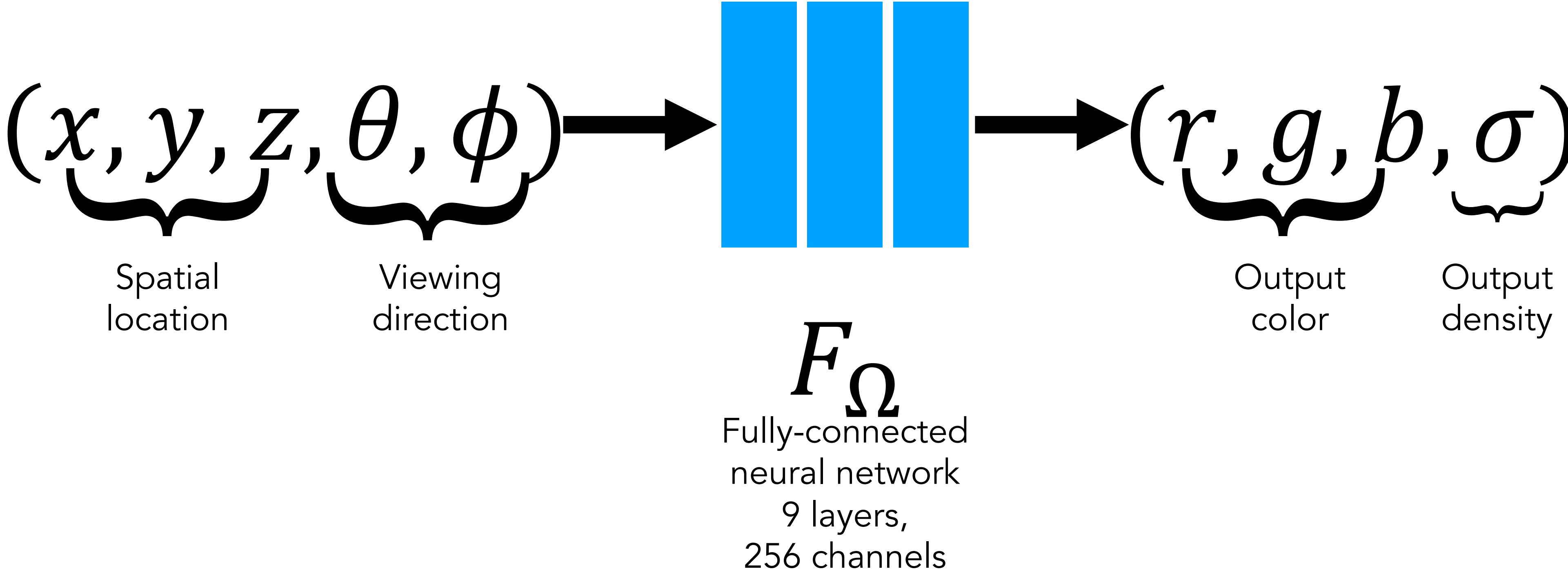
3D-aware Synthesis (part II)

Jun-Yan Zhu

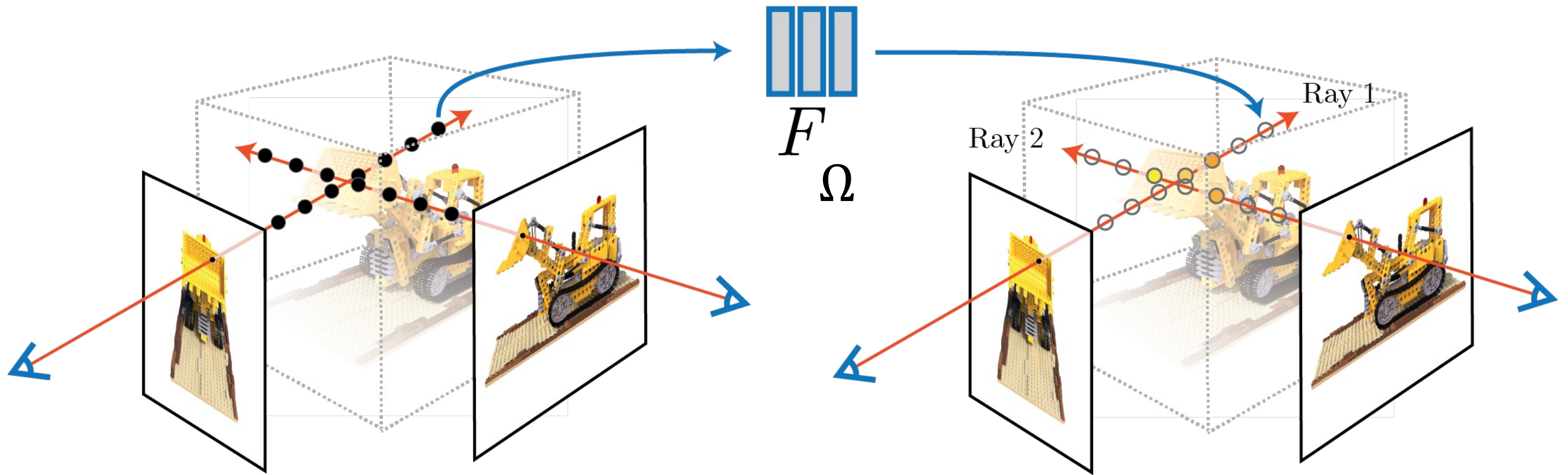
16-726, Spring 2025

NeRF (neural radiance fields):
Neural networks as a *volume* representation,
using volume rendering to do view
synthesis. $(x, y, z, \theta, \phi) \rightarrow \textit{color, opacity}$

Representing a scene as a continuous 5D function

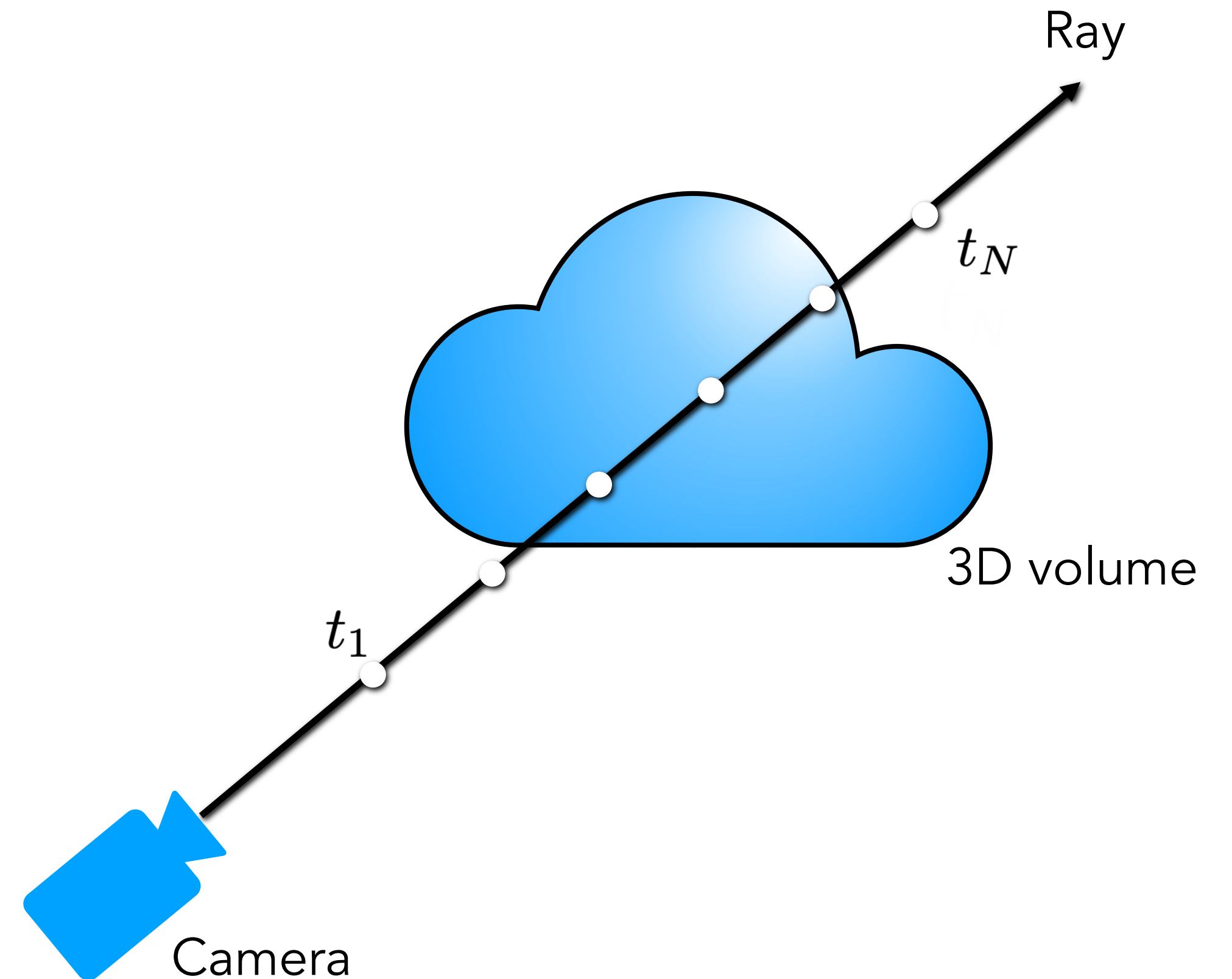


Generate views with traditional volume rendering



Generate views with traditional volume rendering

Rendering model for ray $r(t) = o + td$:

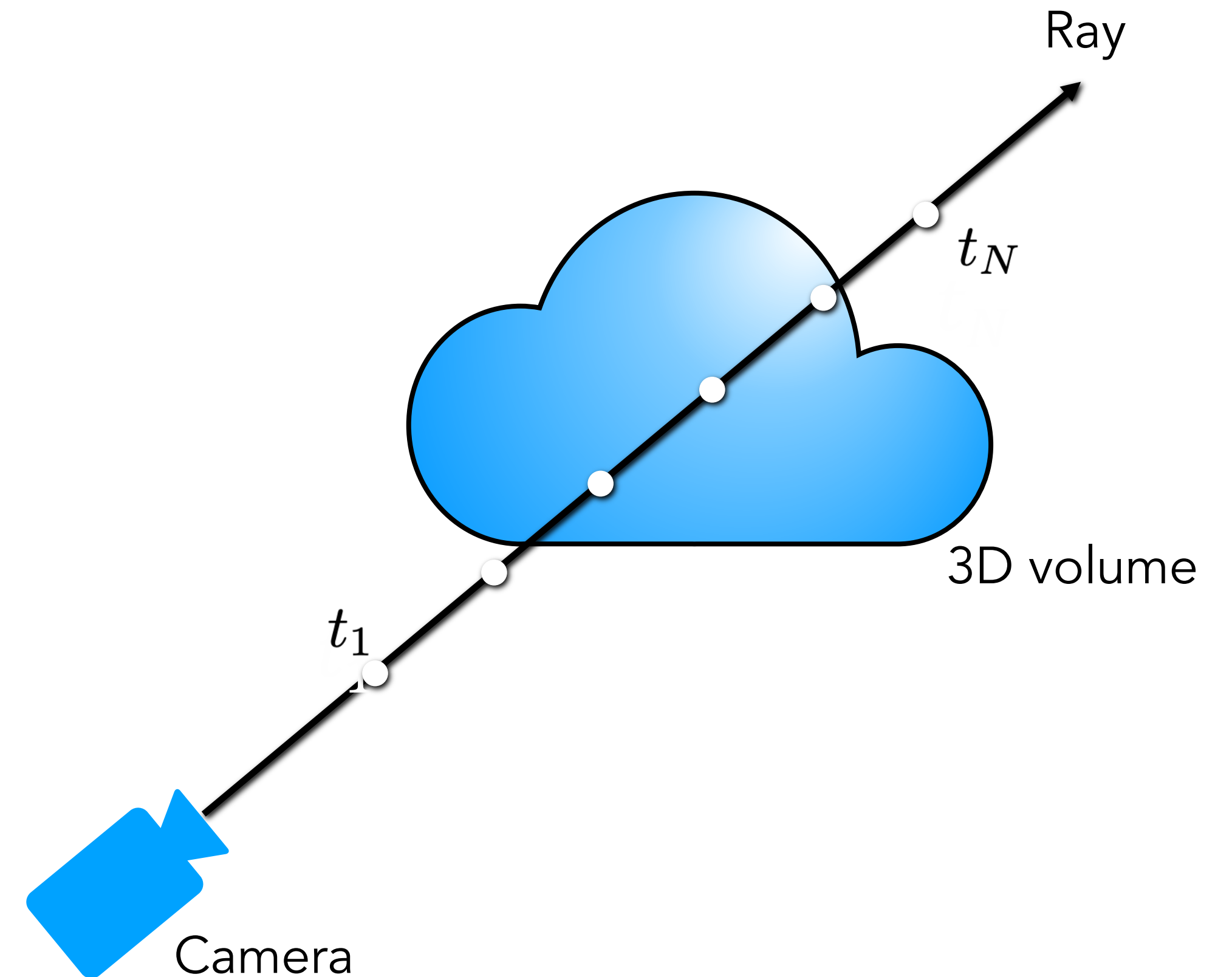


Generate views with traditional volume rendering

Rendering model for ray $r(t) = o + td$:

$$C \approx \sum_{i=1}^N T_i \alpha_i c_i$$

weights colors



Generate views with traditional volume rendering

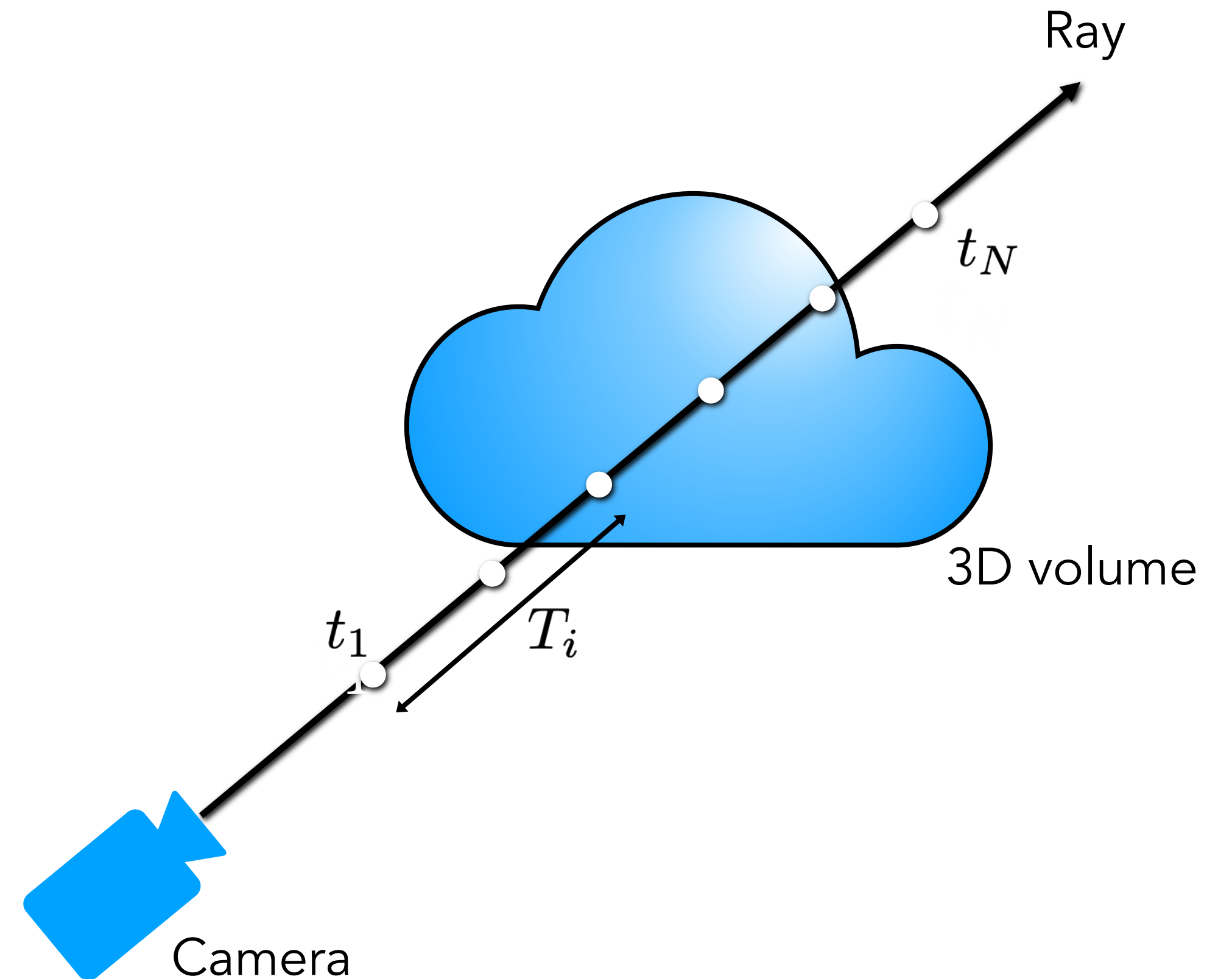
Rendering model for ray $r(t) = o + td$:

$$C \approx \sum_{i=1}^N T_i \alpha_i c_i$$

weights colors

How much light is blocked earlier along ray:

$$T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$$



Generate views with traditional volume rendering

Rendering model for ray $r(t) = o + td$:

$$C \approx \sum_{i=1}^N T_i \alpha_i c_i$$

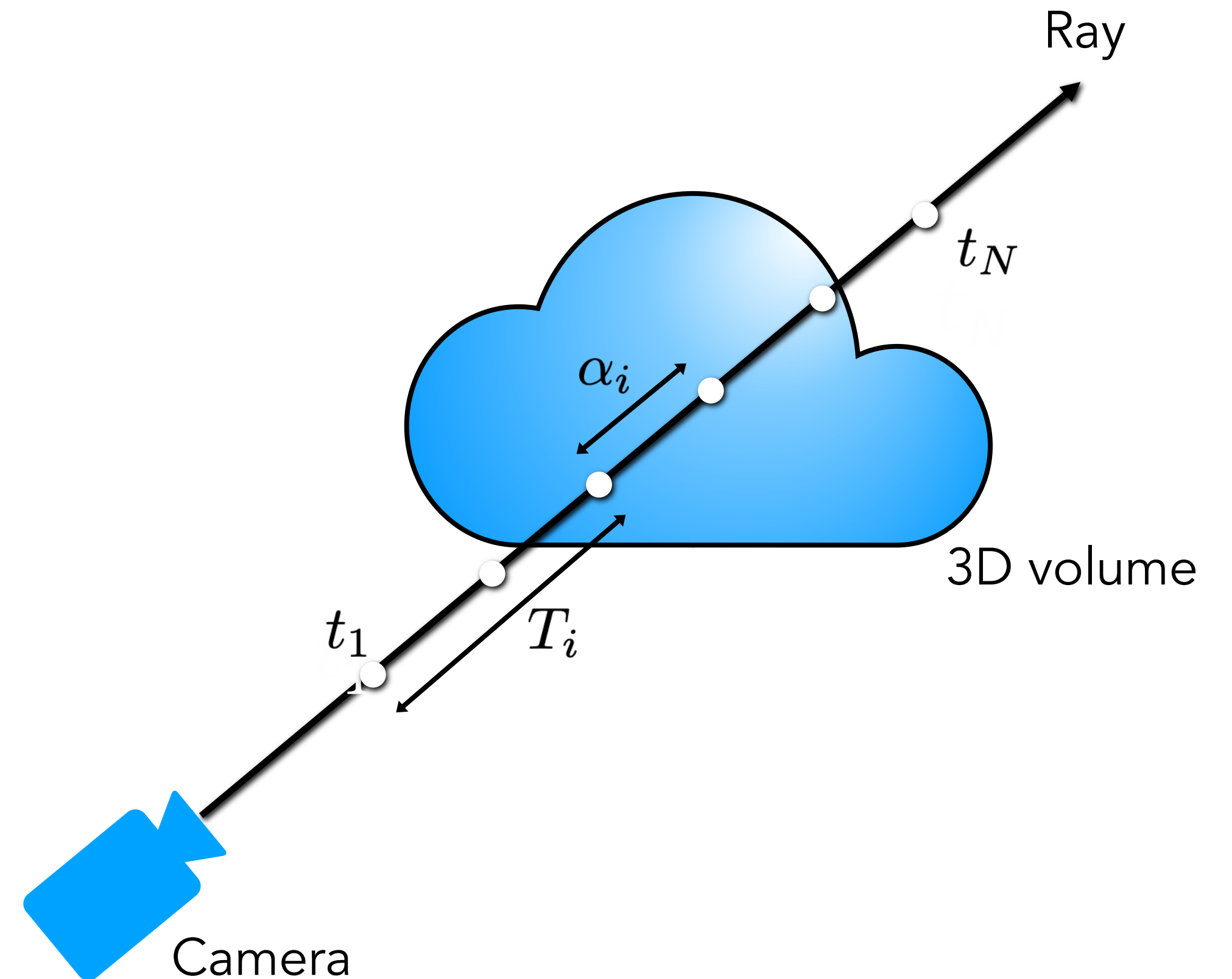
weights colors

How much light is blocked earlier along ray:

$$T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$$

How much light is contributed by ray segment i :

$$\alpha_i = 1 - e^{-\sigma_i \delta t_i}$$



Sigma parametrization for continuous opacity

Rendering model for ray $r(t) = o + td$:

$$C \approx \sum_{i=1}^N T_i \alpha_i c_i$$

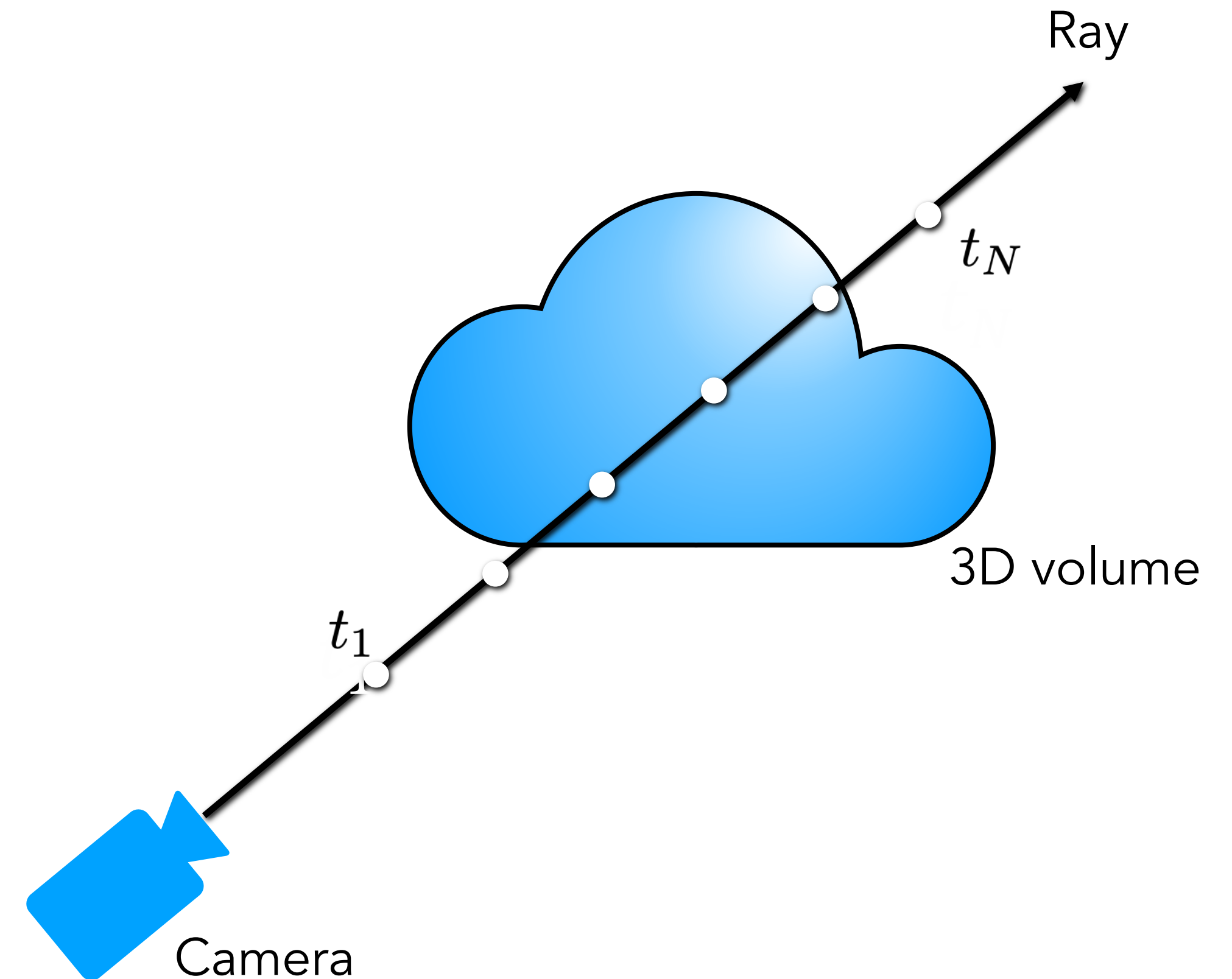
weights colors

How much light is blocked earlier along ray:

$$T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$$

How much light is contributed by ray segment i :

$$\alpha_i = 1 - e^{-\sigma_i \delta t_i}$$



Effective resolution is tied to distance between samples

Rendering model for ray $r(t) = o + td$:

$$C \approx \sum_{i=1}^N T_i \alpha_i c_i$$

weights colors

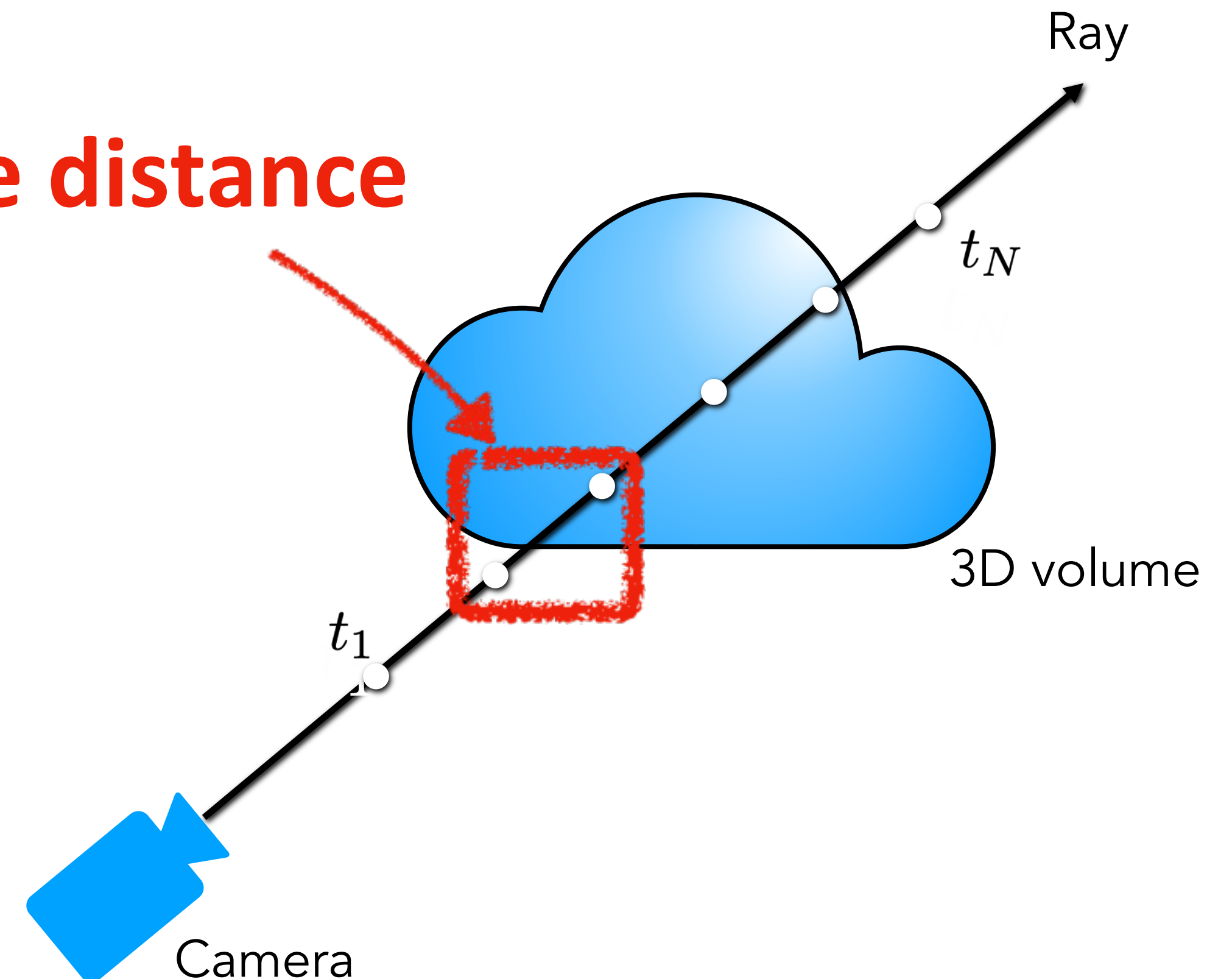
How much light is blocked earlier along ray:

$$T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$$

How much light is contributed by ray segment i :

$$\alpha_i = 1 - e^{-\sigma_i \delta t_i}$$

sample distance



Volume rendering is trivially differentiable

Rendering model for ray $r(t) = o + td$:

$$C \approx \sum_{i=1}^N T_i \alpha_i c_i$$

differentiable w.r.t.

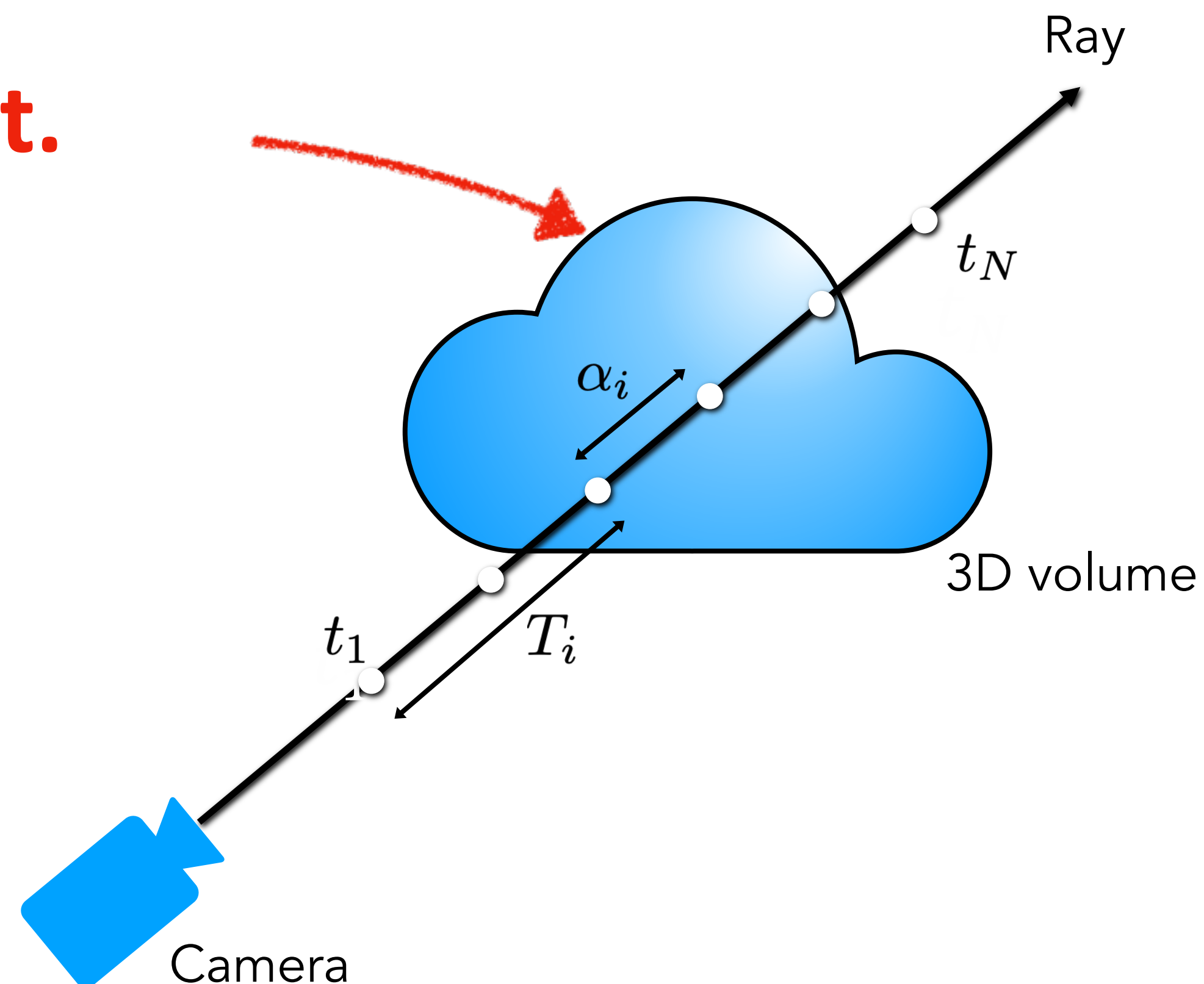
weights → T_i
colors → c_i

How much light is blocked earlier along ray:

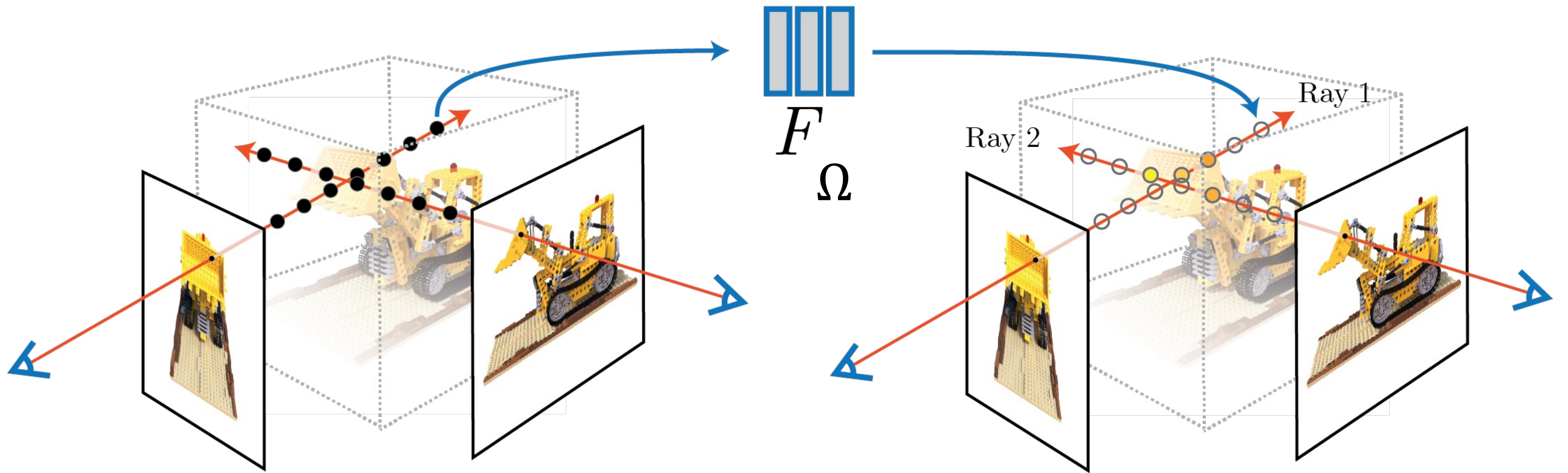
$$T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$$

How much light is contributed by ray segment i :

$$\alpha_i = 1 - e^{-\sigma_i \delta t_i}$$



Optimize with gradient descent on rendering loss

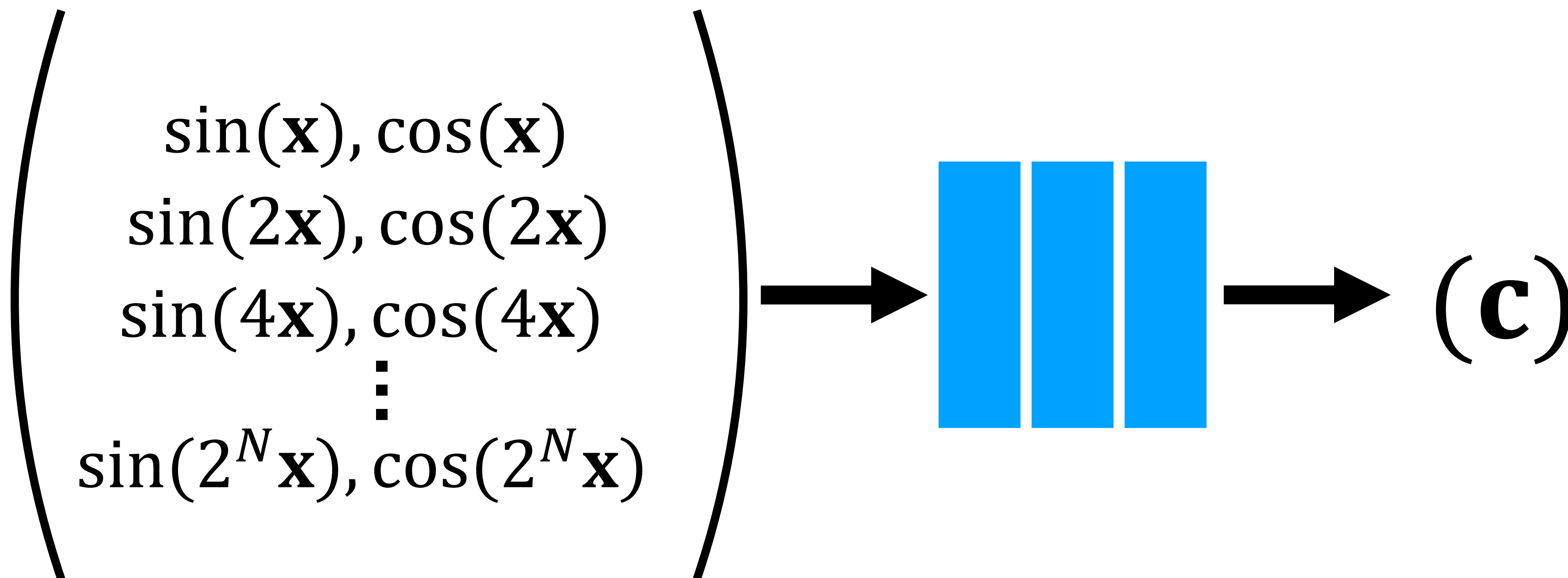


$$\min_{\Omega} \sum_i \left\| \text{render}^{(i)}(F_{\Omega}) - I_{\text{gt}}^{(i)} \right\|^2$$

Training network to reproduce all input views of the scene



Positional encoding: high frequency embedding of input coordinates



Simple trick enables network to memorize images

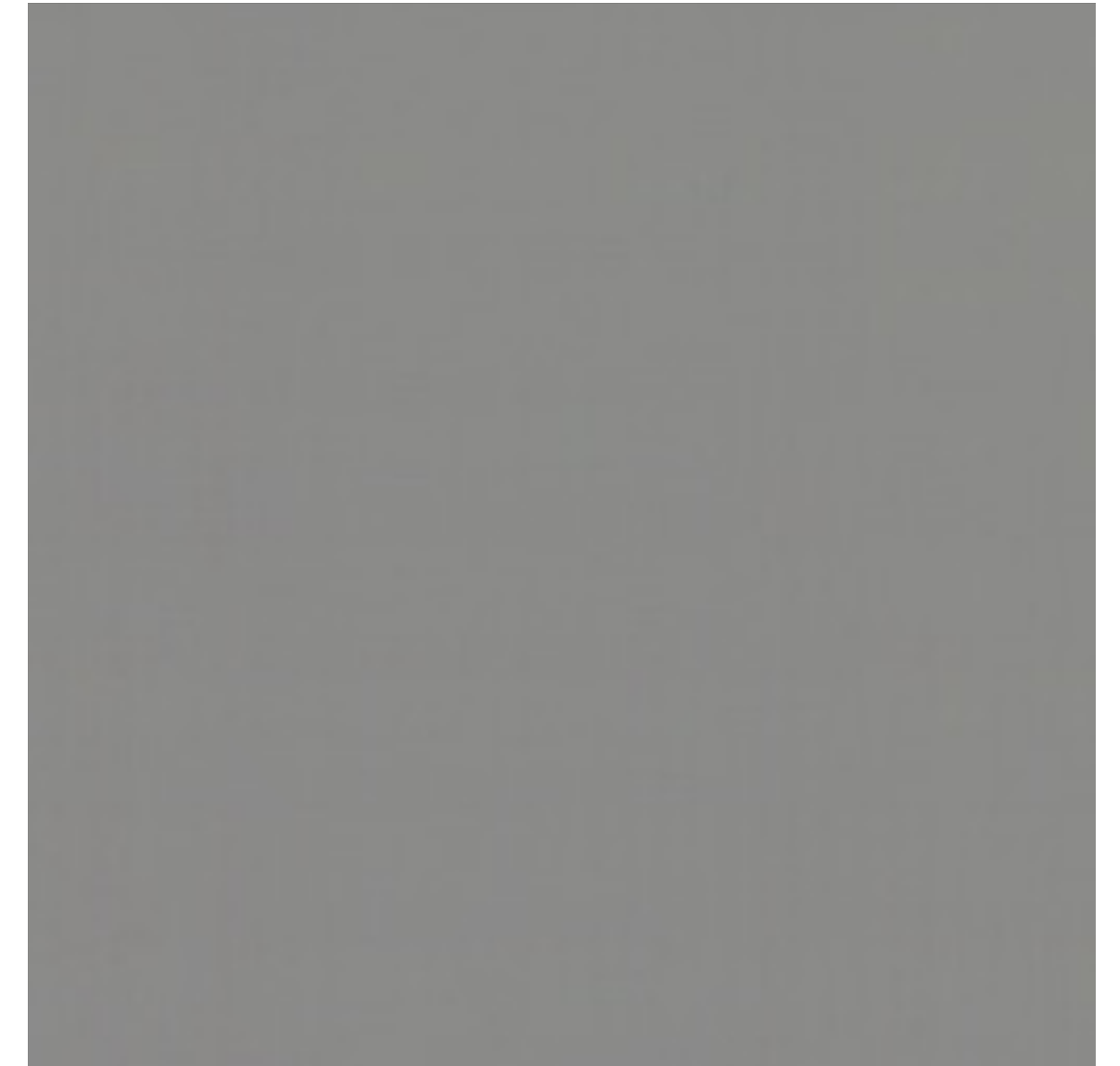
Ground truth image



Standard fully-connected net



With "embedding"



Positional encoding also directly improves our scene representation!



NeRF (Naive)



NeRF (with positional encoding)

Implementation Details

Camera Locations and Poses

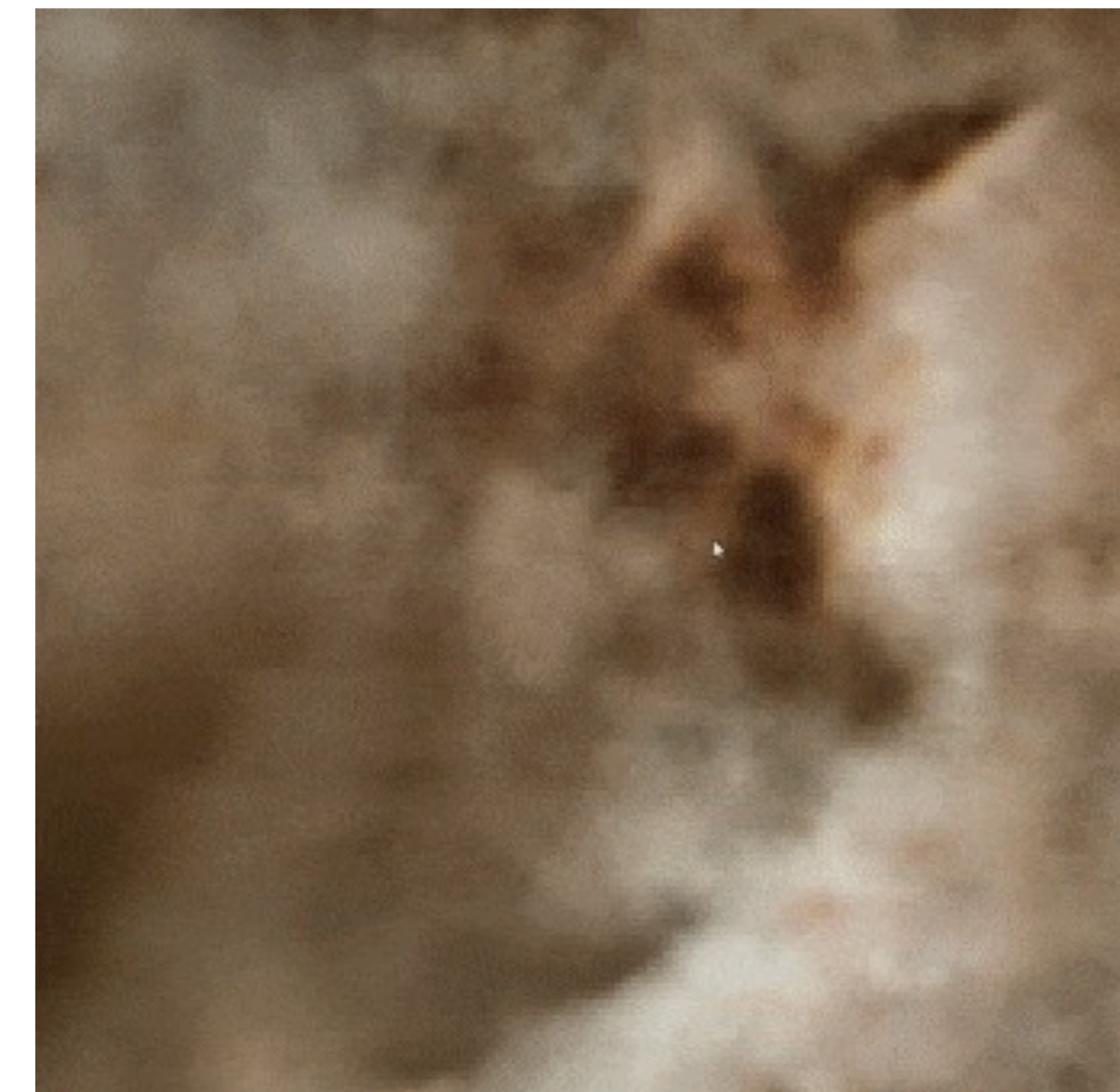
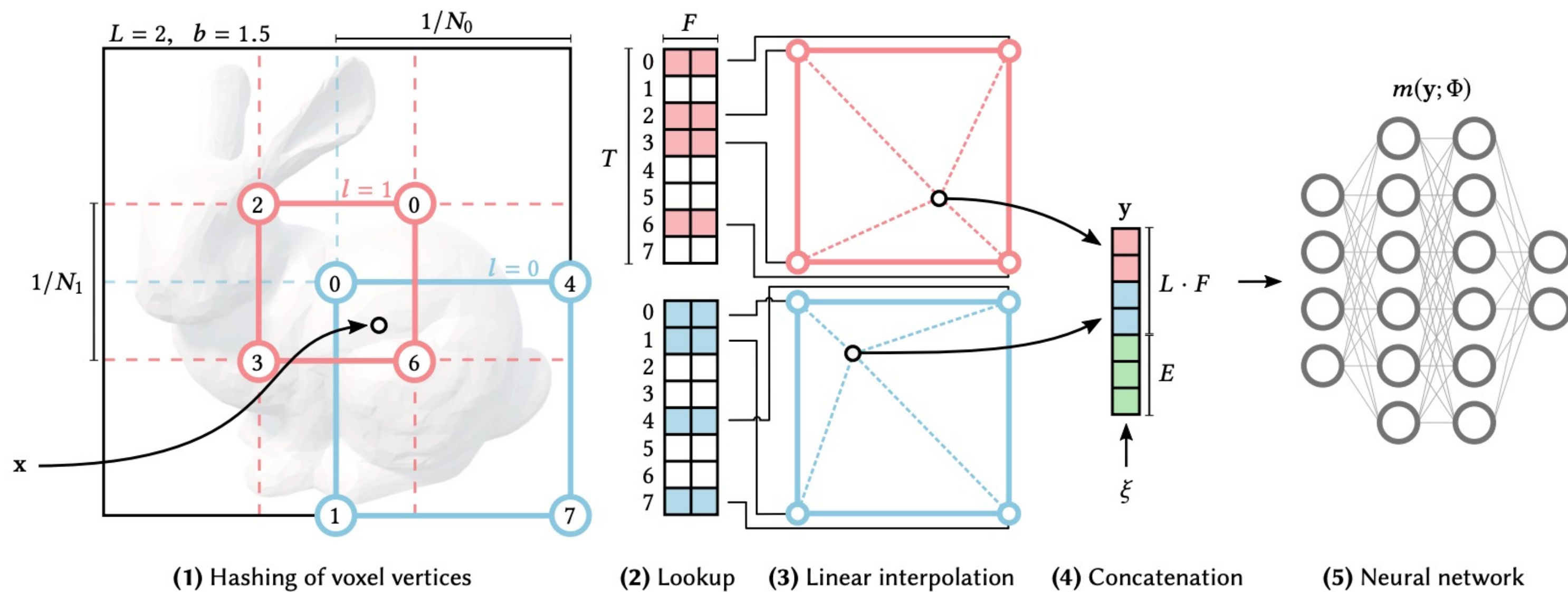
- Use Structure from Motion (e.g., [COLMAP](#)) to initialize camera poses
- Incorrect camera poses lead to bad results
- Joint optimization of camera poses and scene presentation.



Implementation Details

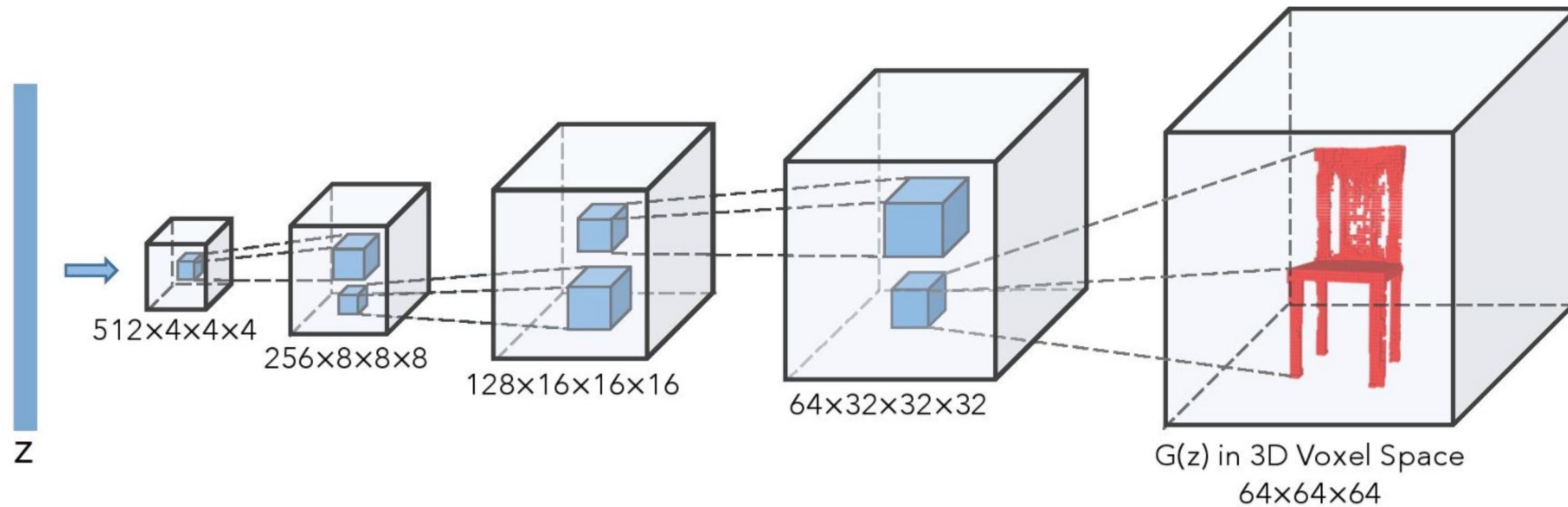
Training and inference speed:

- Original NeRF is quite slow.
- Faster training and inference is an active research topic.
- Optimized CUDA kernel for small MLP network (10x faster)
- Efficient data structure: multi-resolution hashing (10+ faster)

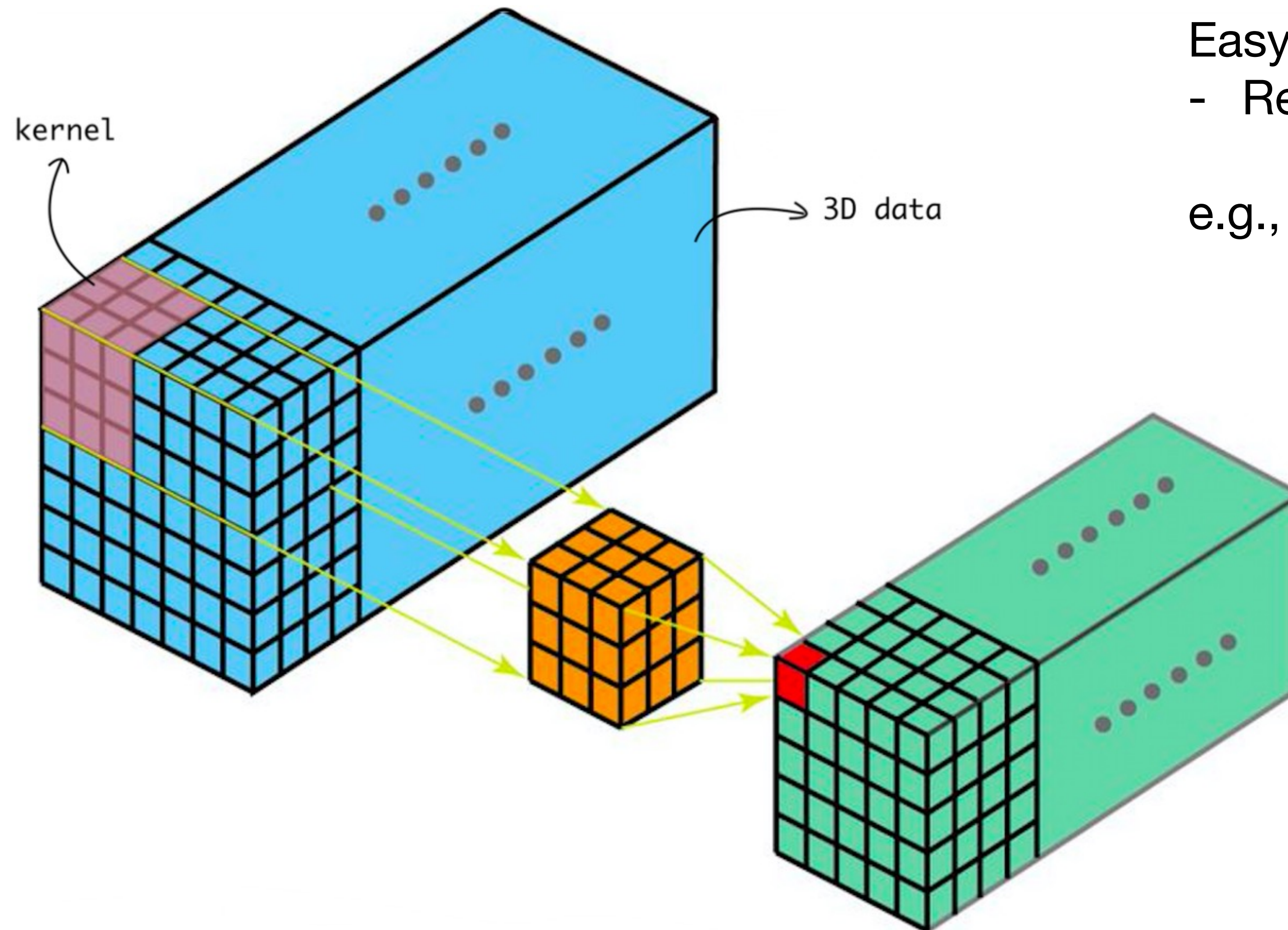


Toward 3D-aware Generative Models

3D Generative Adversarial Networks



3D Convolutional Layers



Easy to implement:

- Replace 2D by 3D in your code

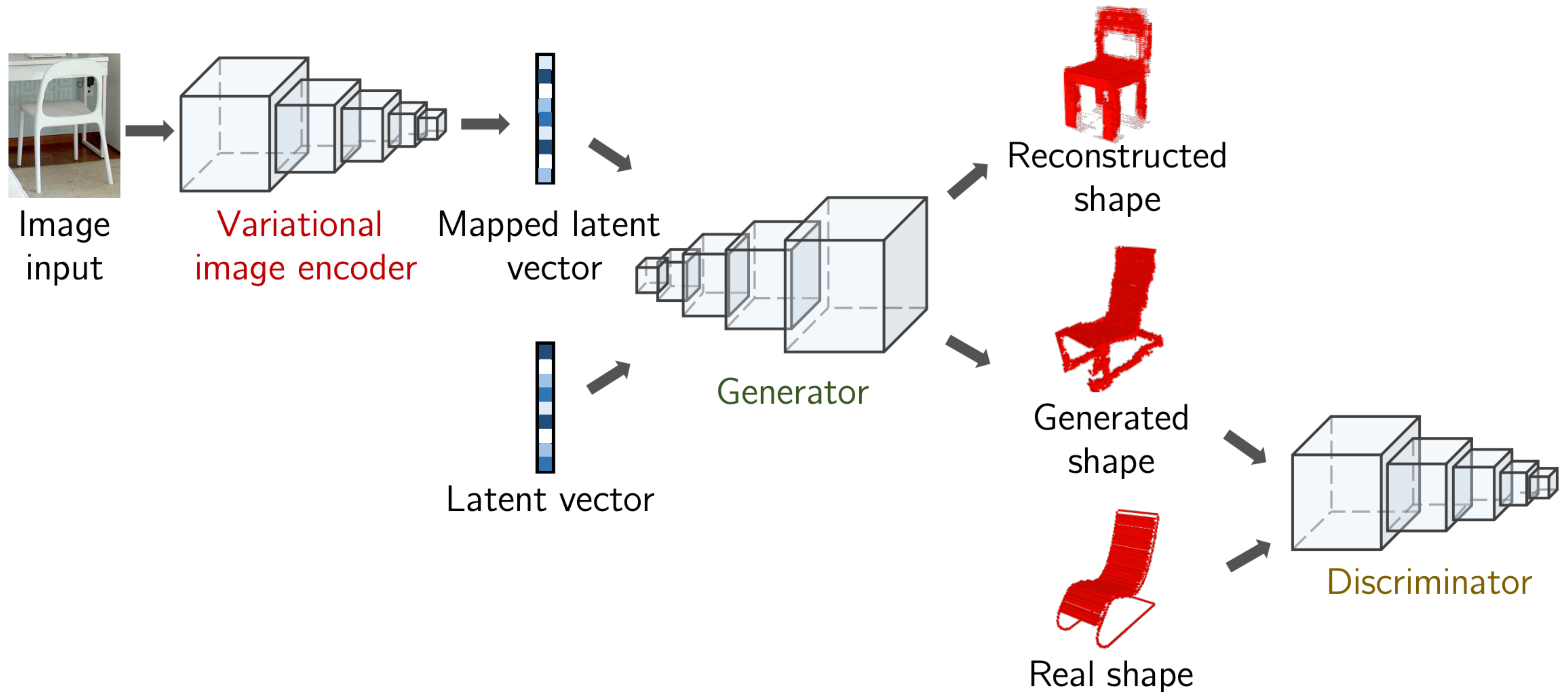
e.g., Conv2D -> Conv3D

ConvTranspose2d->ConvTranspose3d

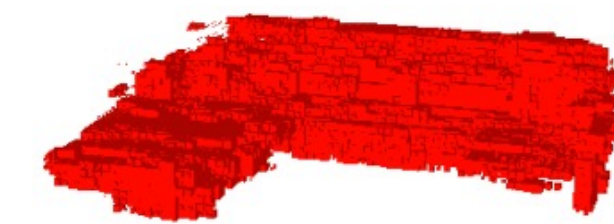
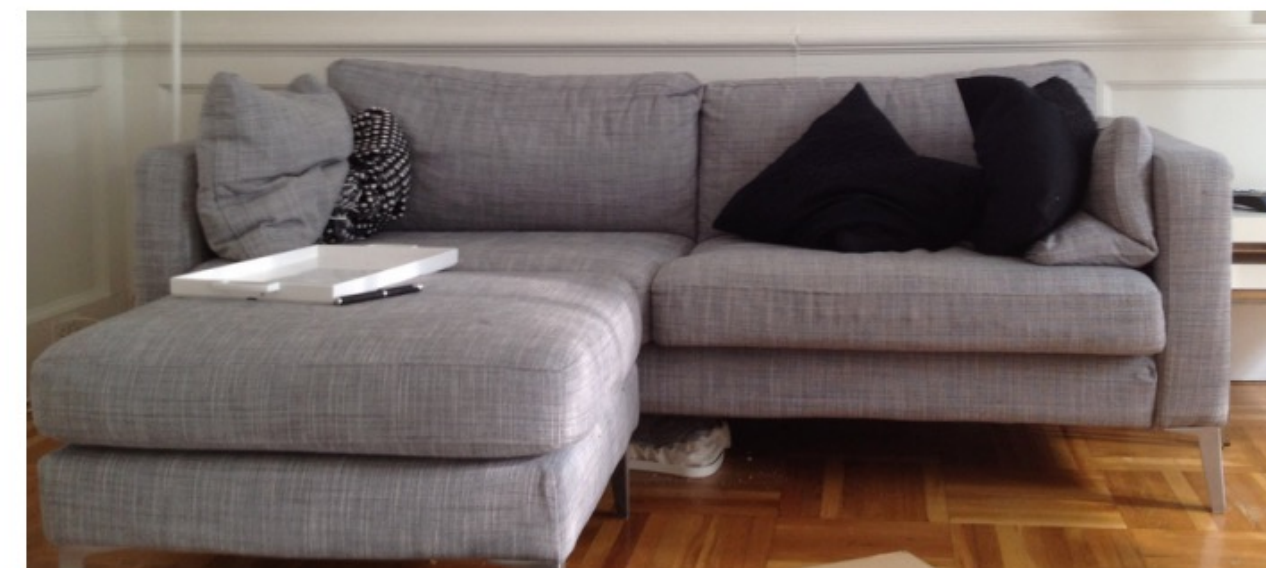
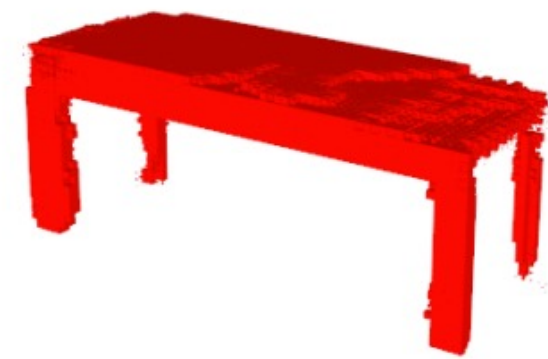
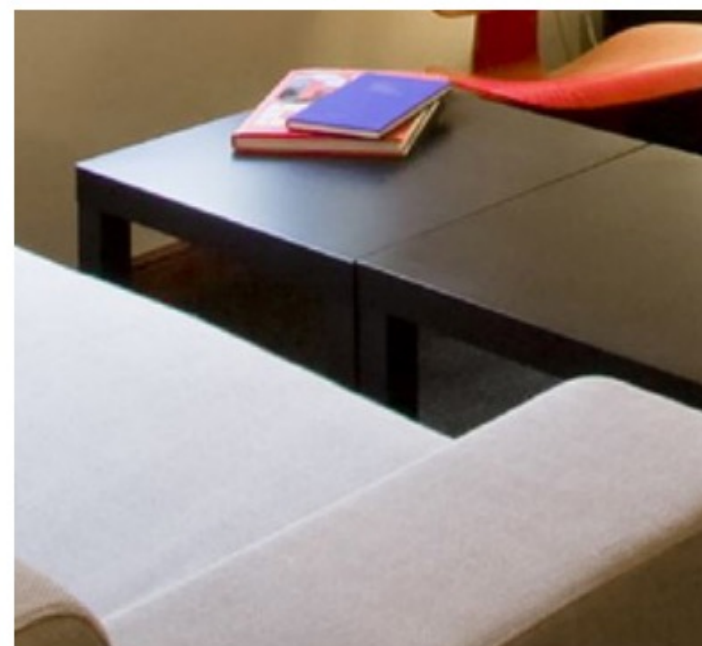
MaxPool2d -> MaxPool3d

```
CLASS torch.nn.Conv3d(in_channels, out_channels, kernel_size, stride=1, padding=0,  
dilation=1, groups=1, bias=True, padding_mode='zeros', device=None, dtype=None) [SOURCE]
```

3D Generative Adversarial Networks



3D Generative Adversarial Networks



Input
image

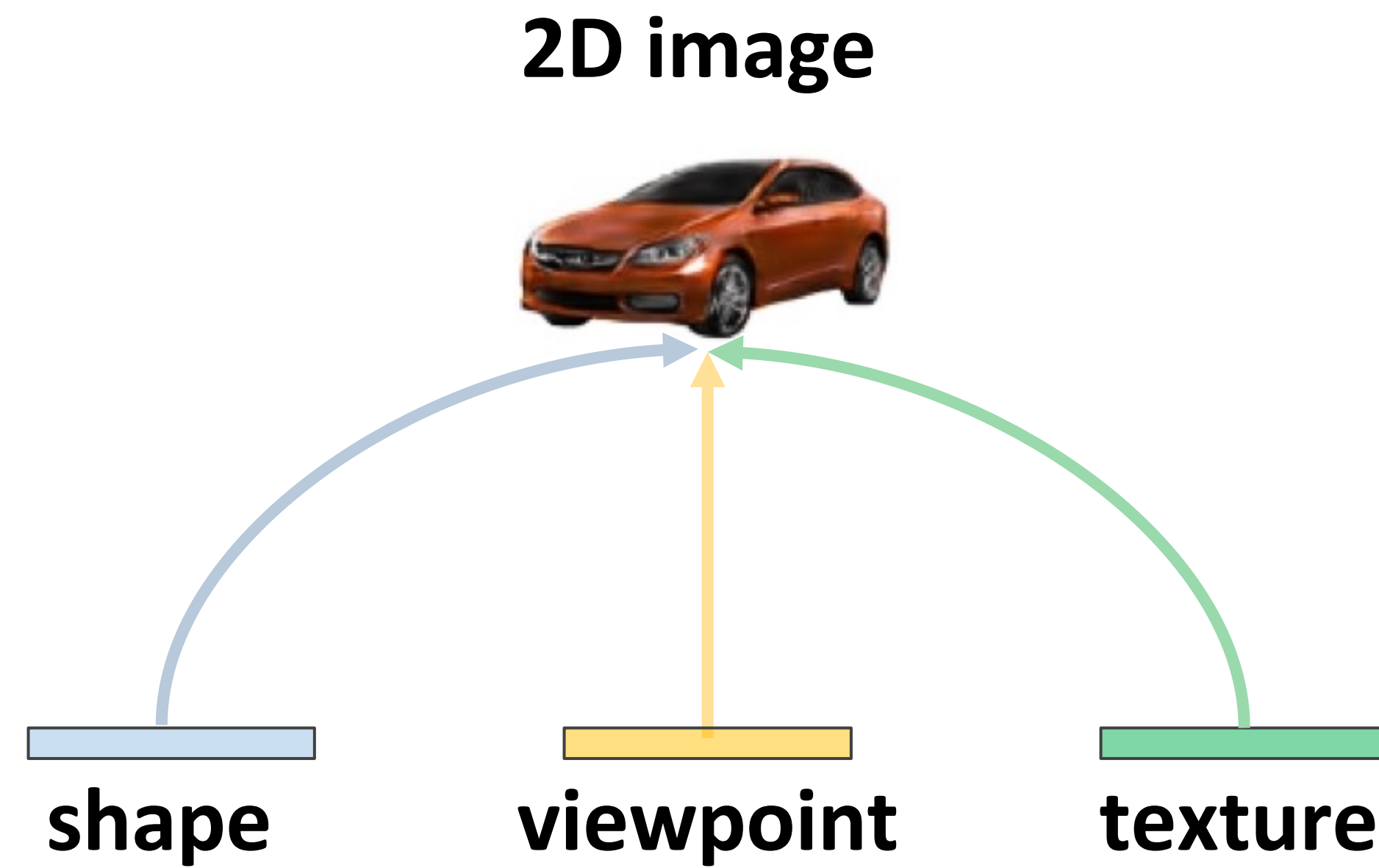
Reconstructed
3D shape

Input
image

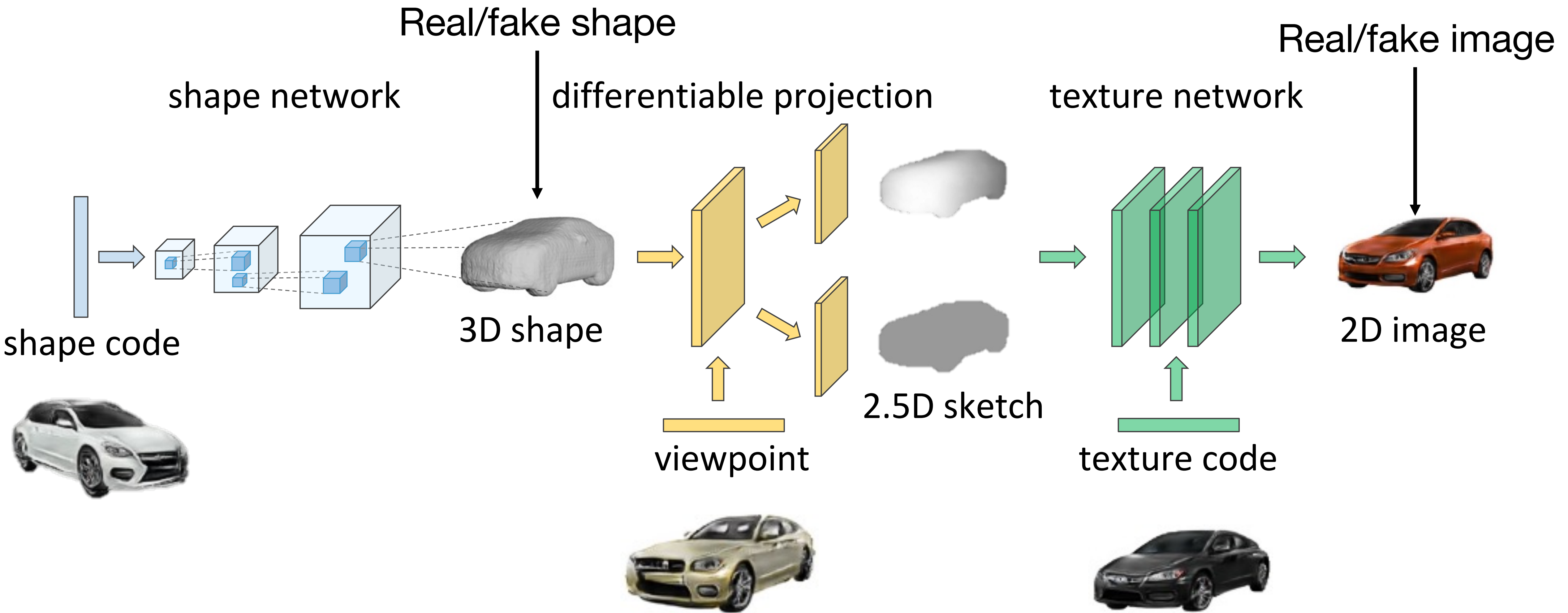
Reconstructed
3D shape

How to add Color and Texture?

Learning 3D Disentanglement



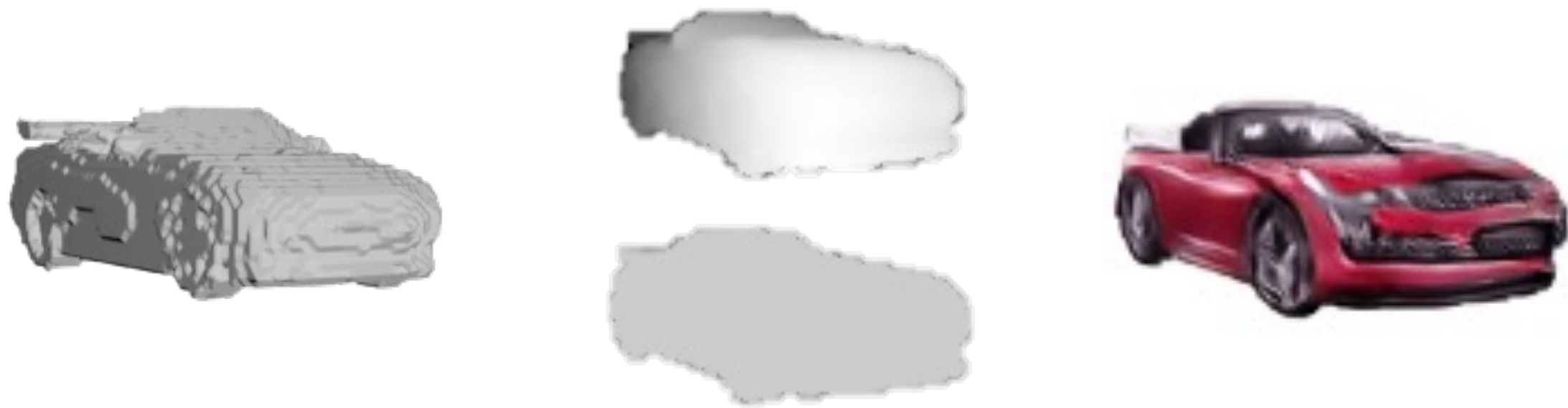
Learning 3D Disentanglement



Learning 3D Disentanglement



samples from 2D GANs



our 3D, 2.5D, and 2D output

viewpoint



shape

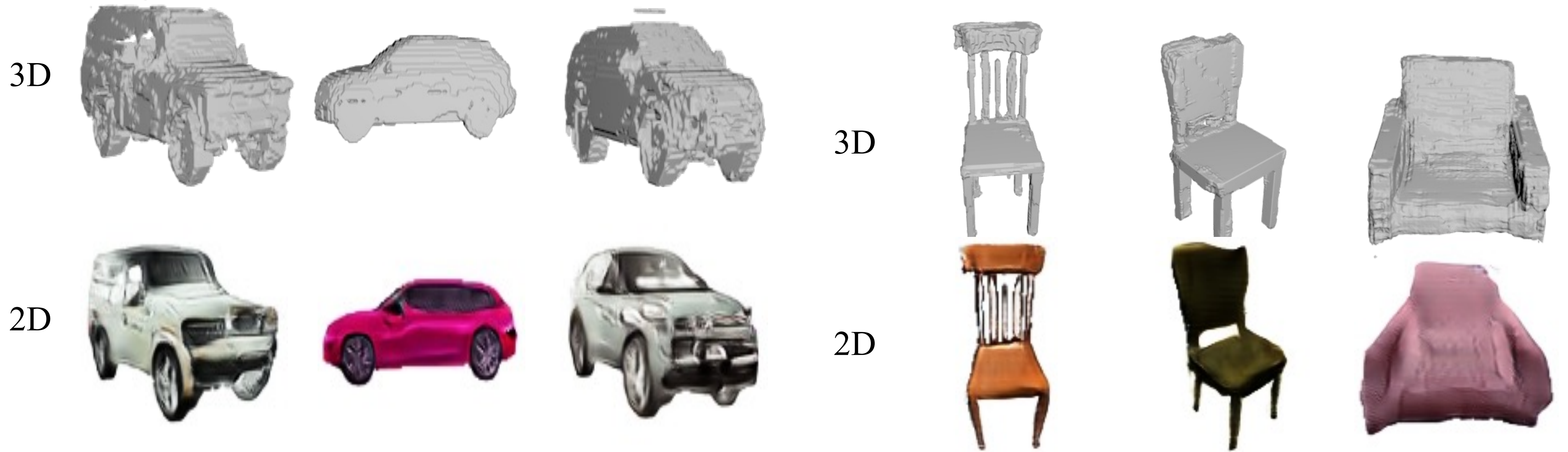


texture

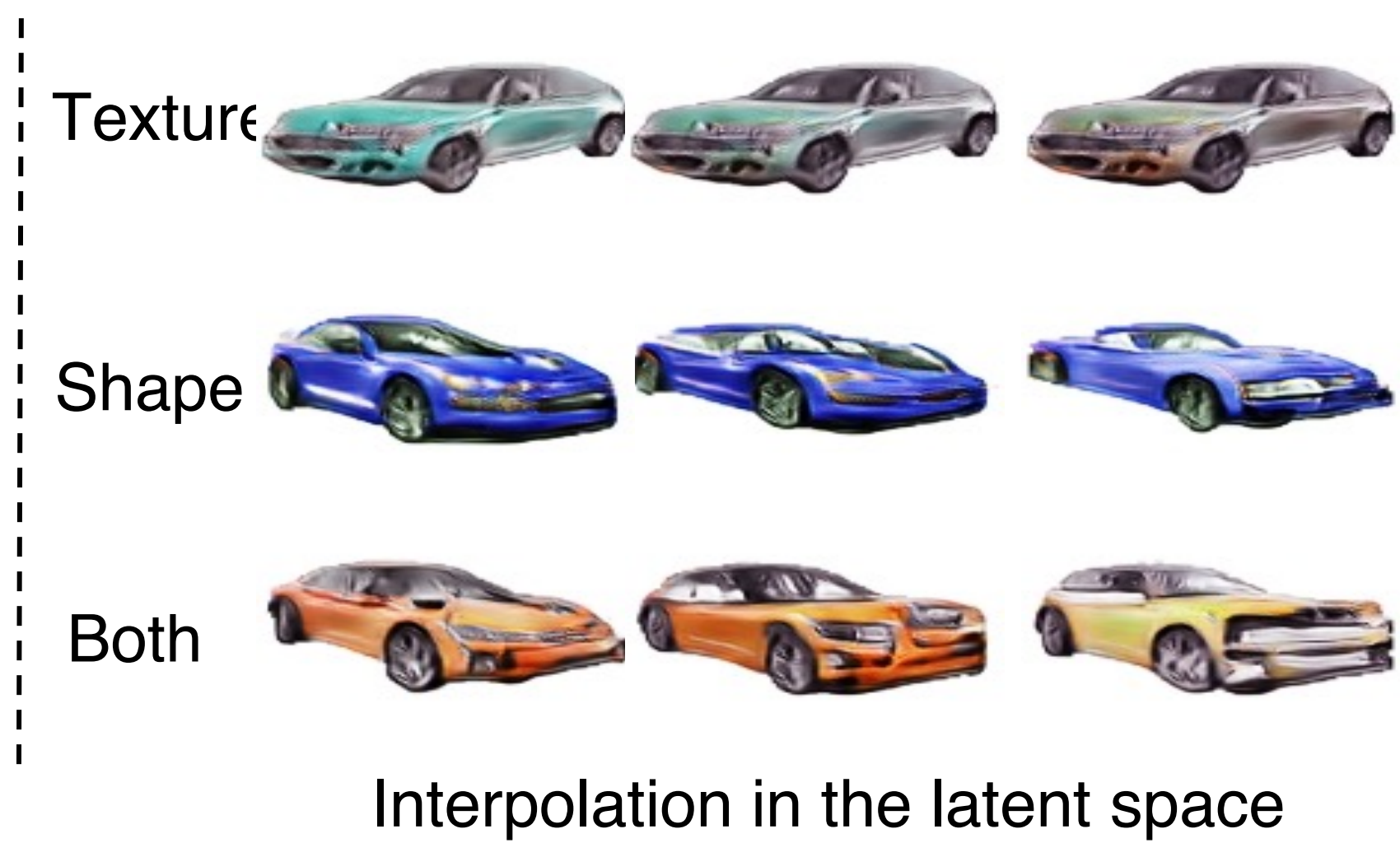
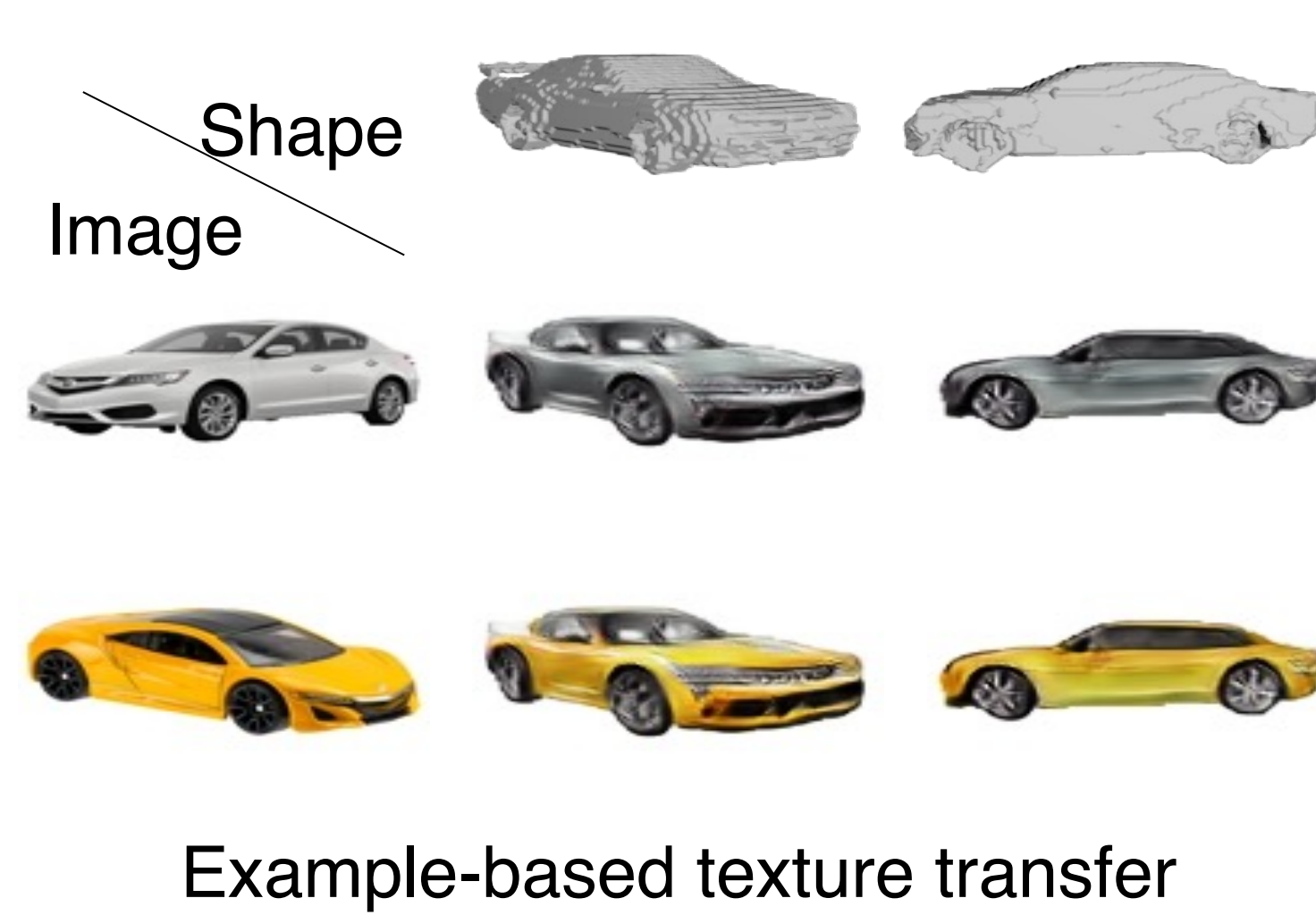
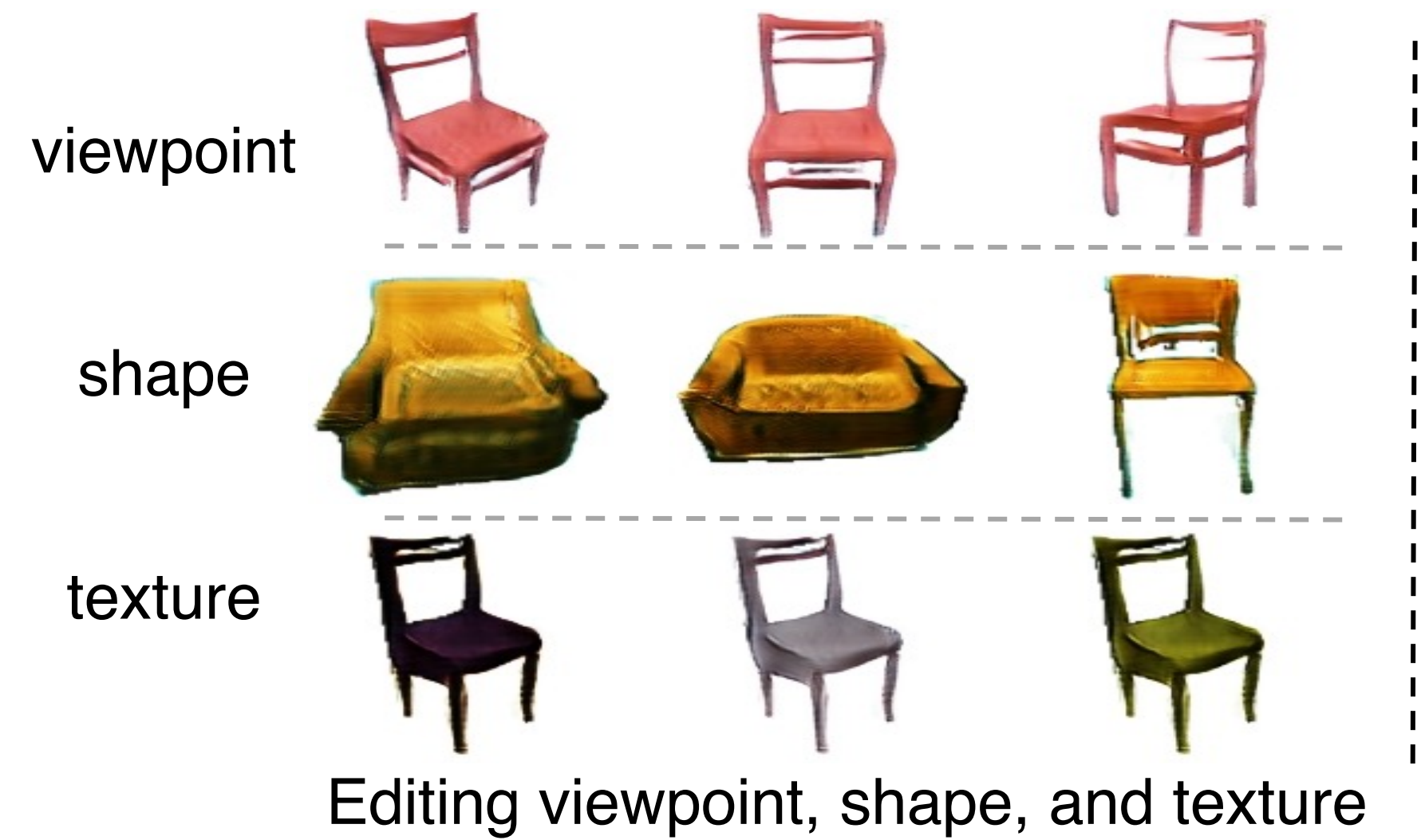


3D disentanglement

Learning 3D Disentanglement



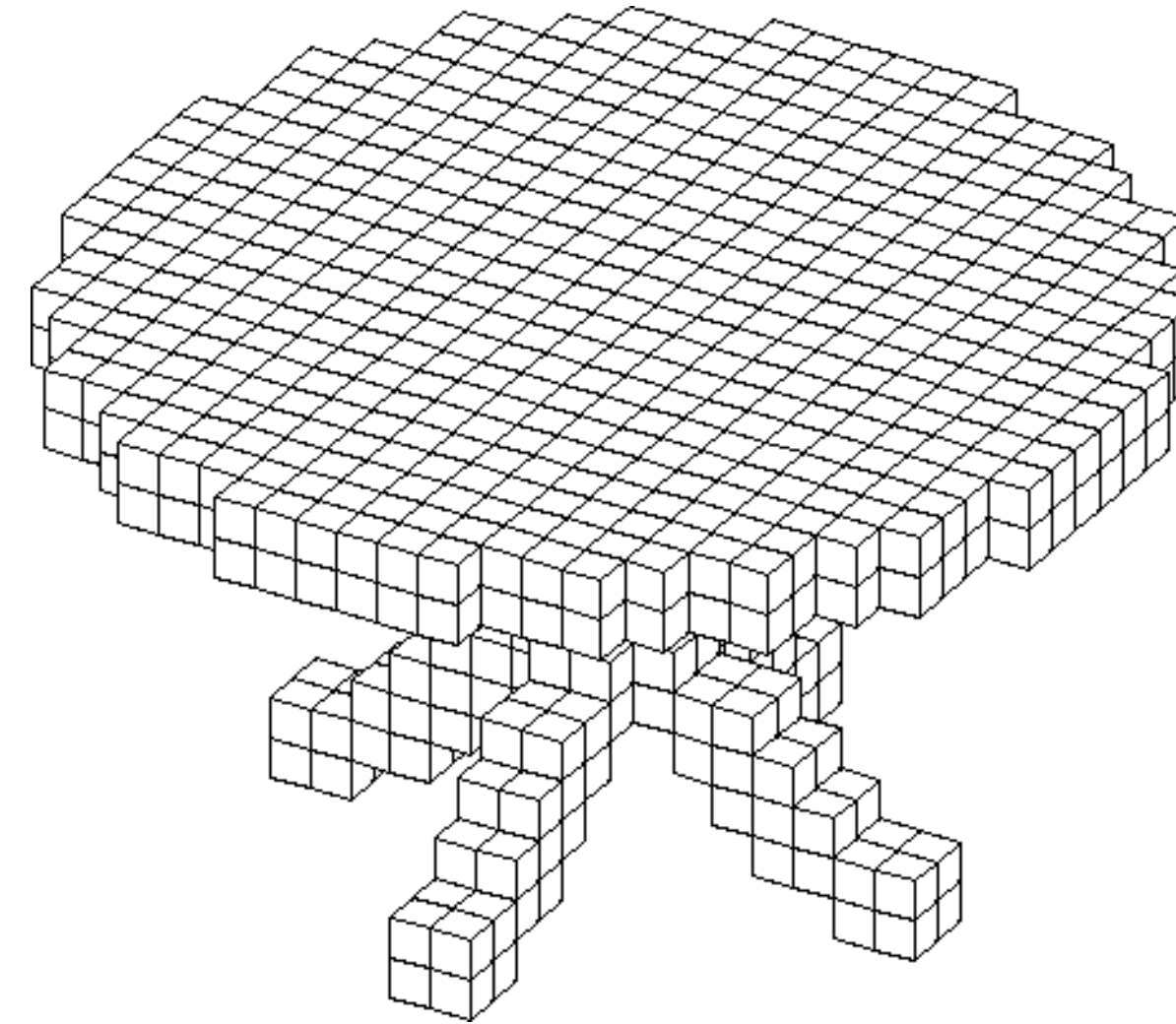
Learning 3D Disentanglement



Limitations:

1. Voxel representation is expensive.
2. Requires ground truth 3D data.

Volumetric 3D



Each grid cell stores information (e.g., occupancy, color)

Very general but memory-intensive

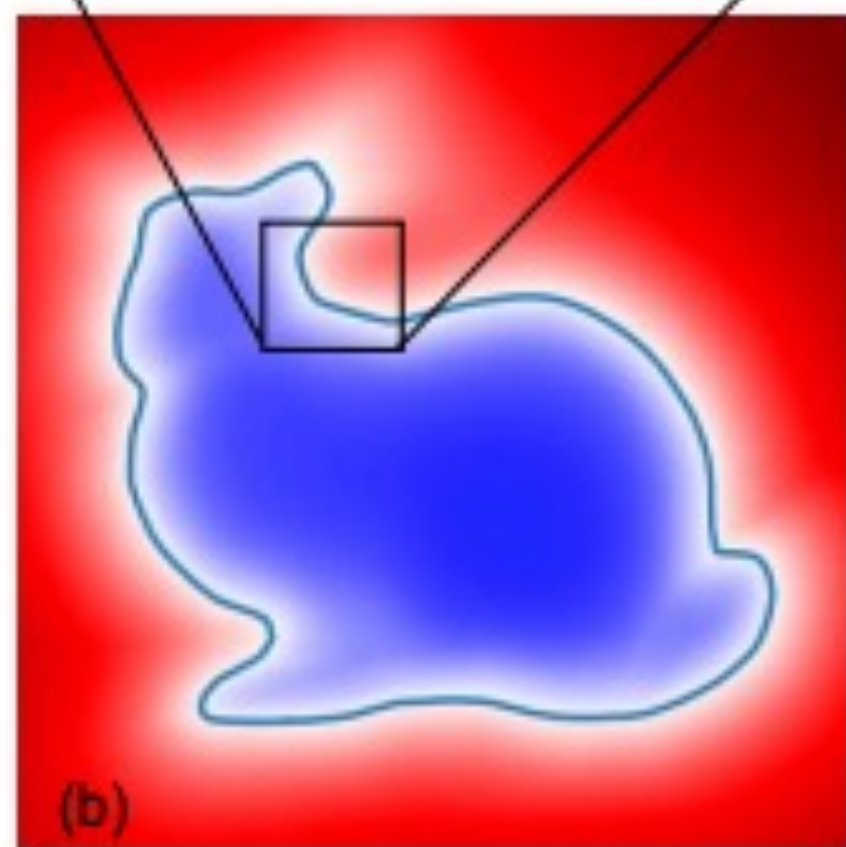
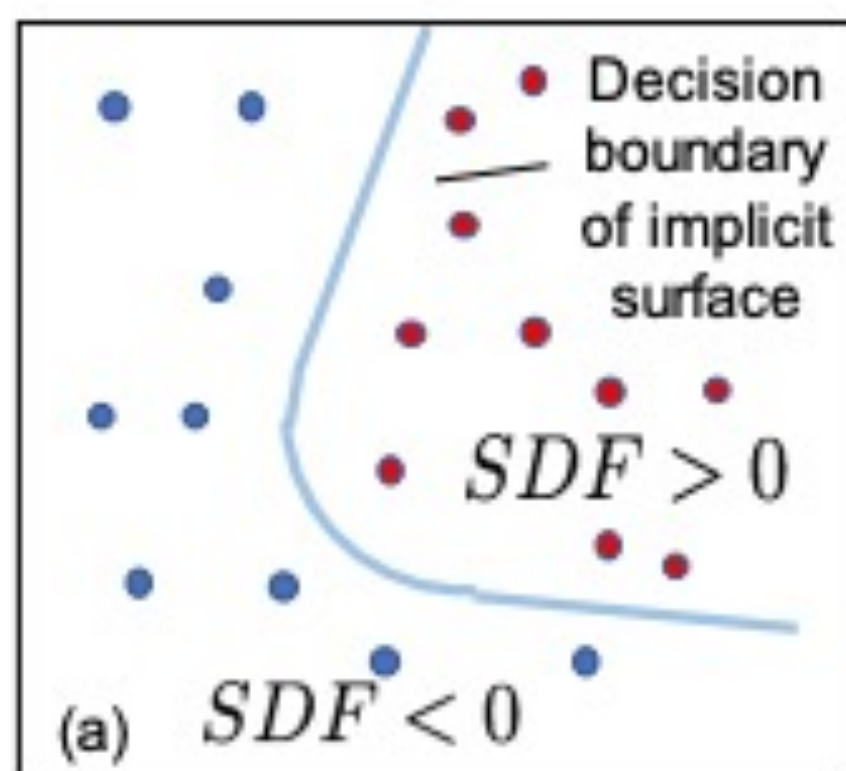
$256 \times 256 \times 256 \rightarrow 1024 \times 1024 \times 1024$

Cannot even fit a single training data to GPU

Improvements:

1. Using implicit representation (network-based)

Signed Distance Function (SDF)



Explicit function:

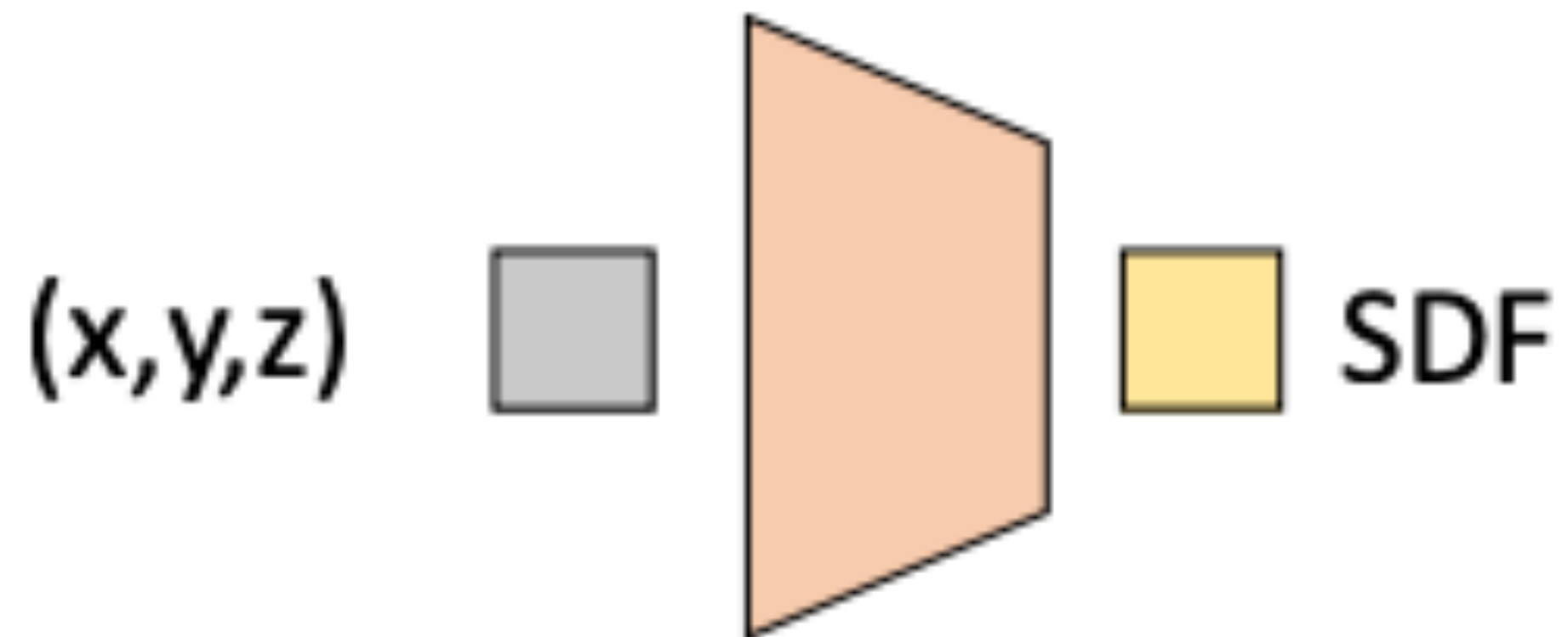
$$y = 2x. (y = f(x))$$

Implicit function:

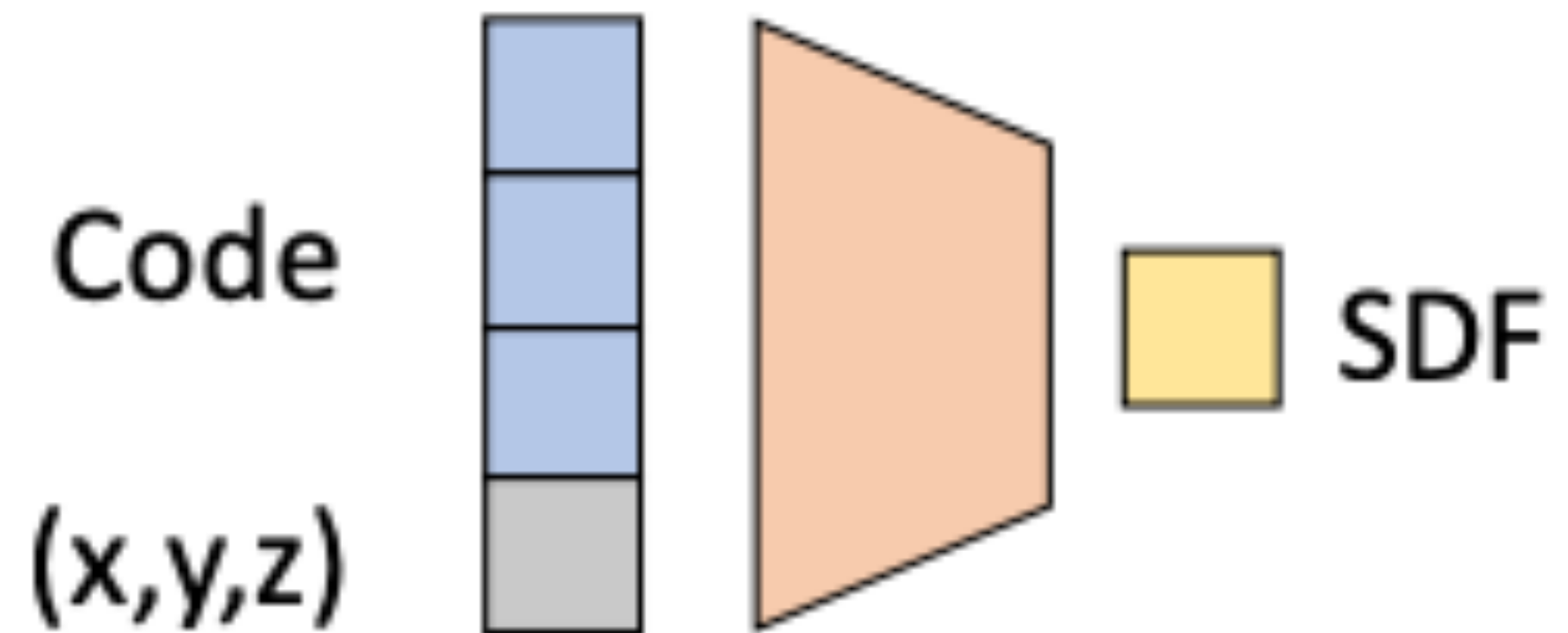
$$2y - 4x = 0, F(x, y) = 0$$

A set of zeros of a function of two variables.

Deep SDF

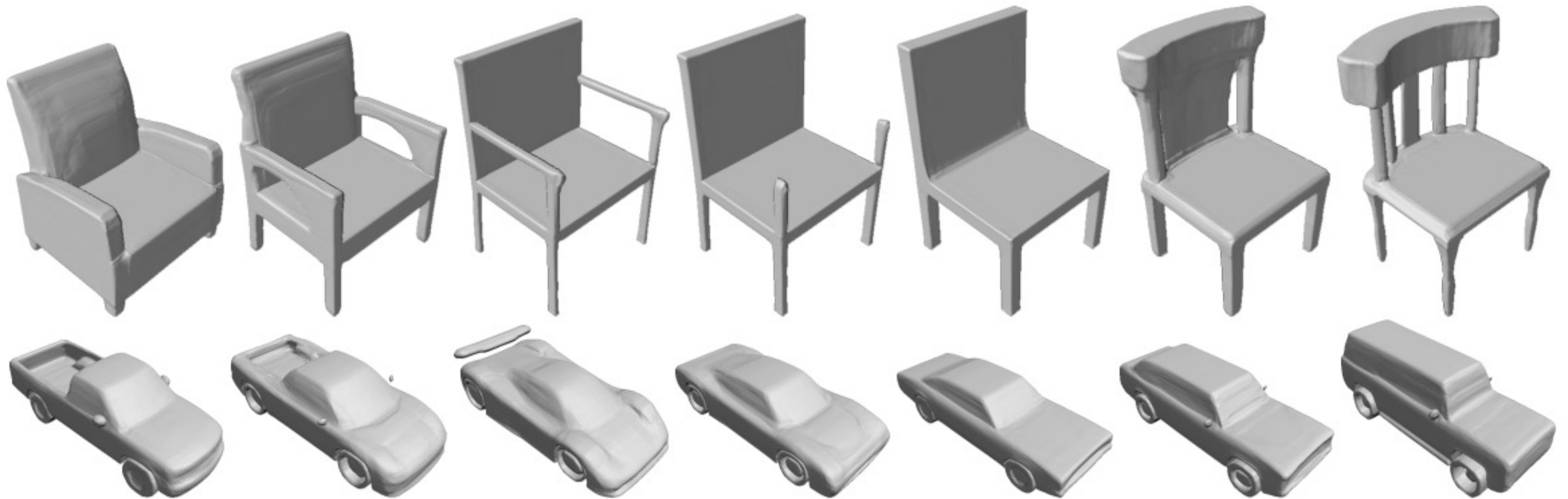


(a) Single Shape DeepSDF

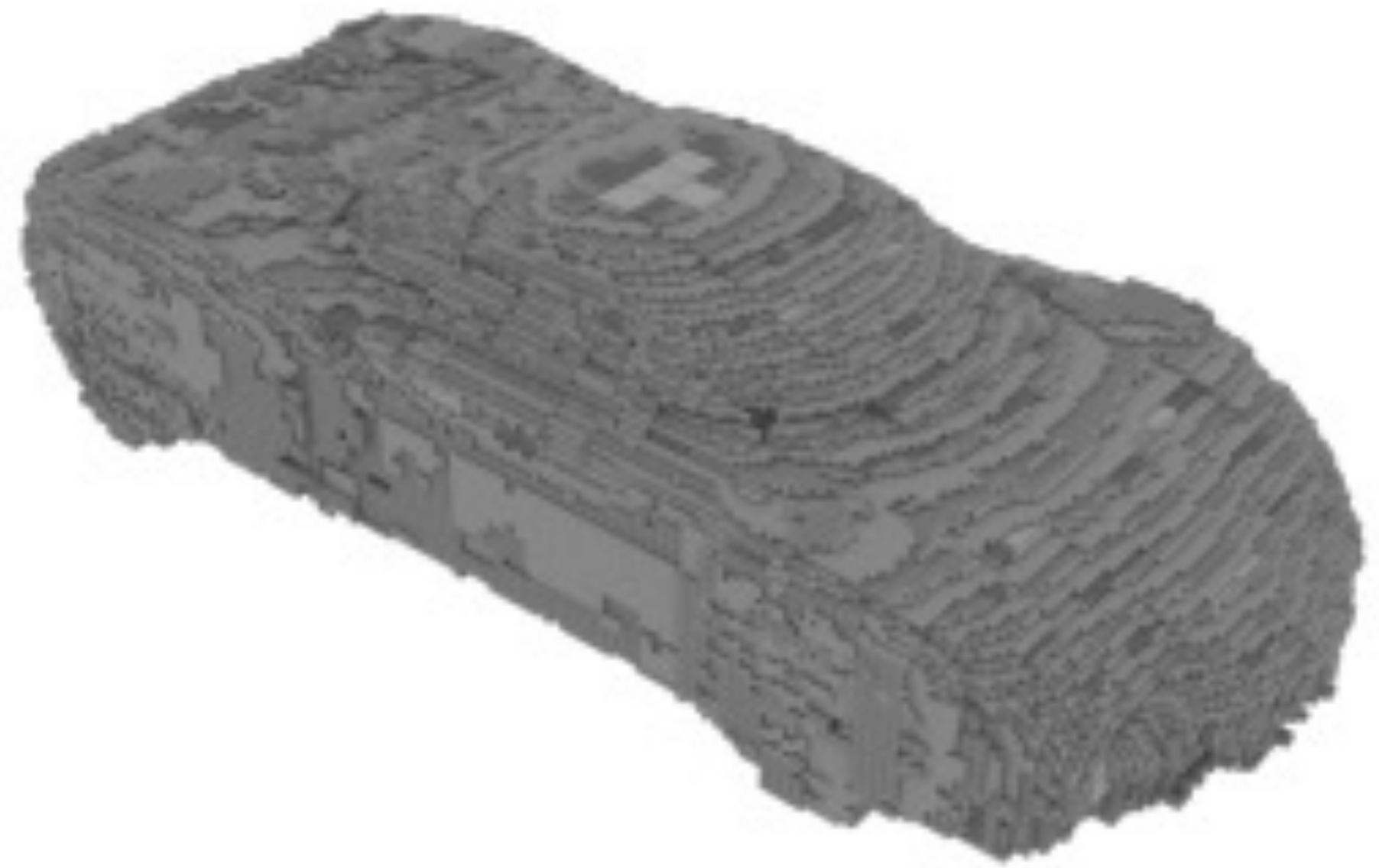


(b) Coded Shape DeepSDF

Deep SDF



Deep SDF

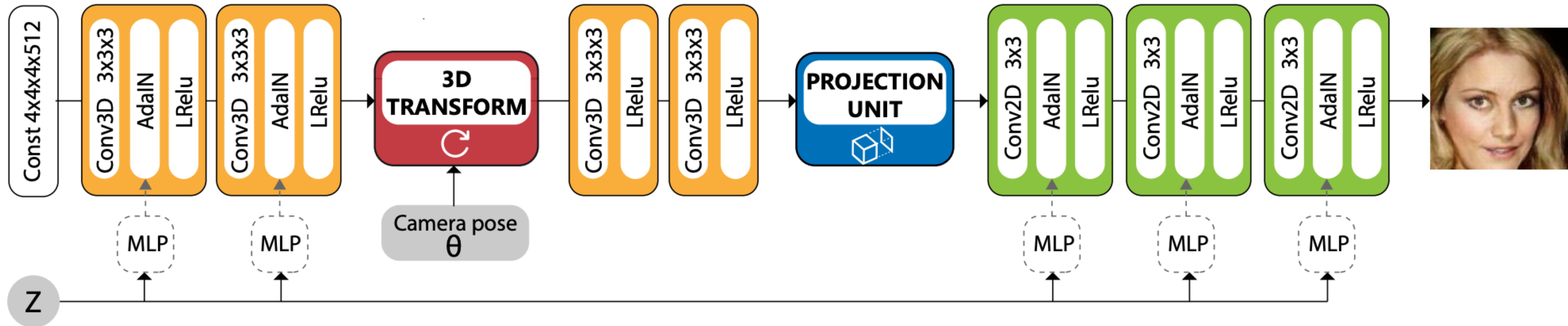


DeepSDF preserve details and render visually pleasing results compared to voxel-based methods.

Improvements:

1. Using implicit representation (network-based)
2. Learning from image collections

HoloGAN

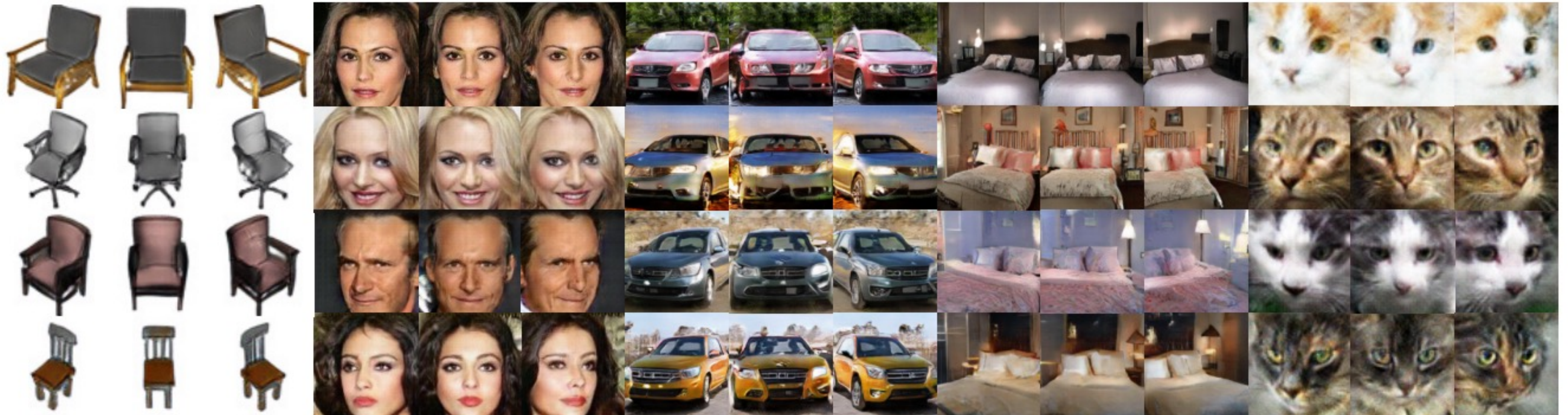


Representation: 3D feature representation

Training: Adversarial loss + latent code reconstruction

Modulation: AdaIN

HoloGAN



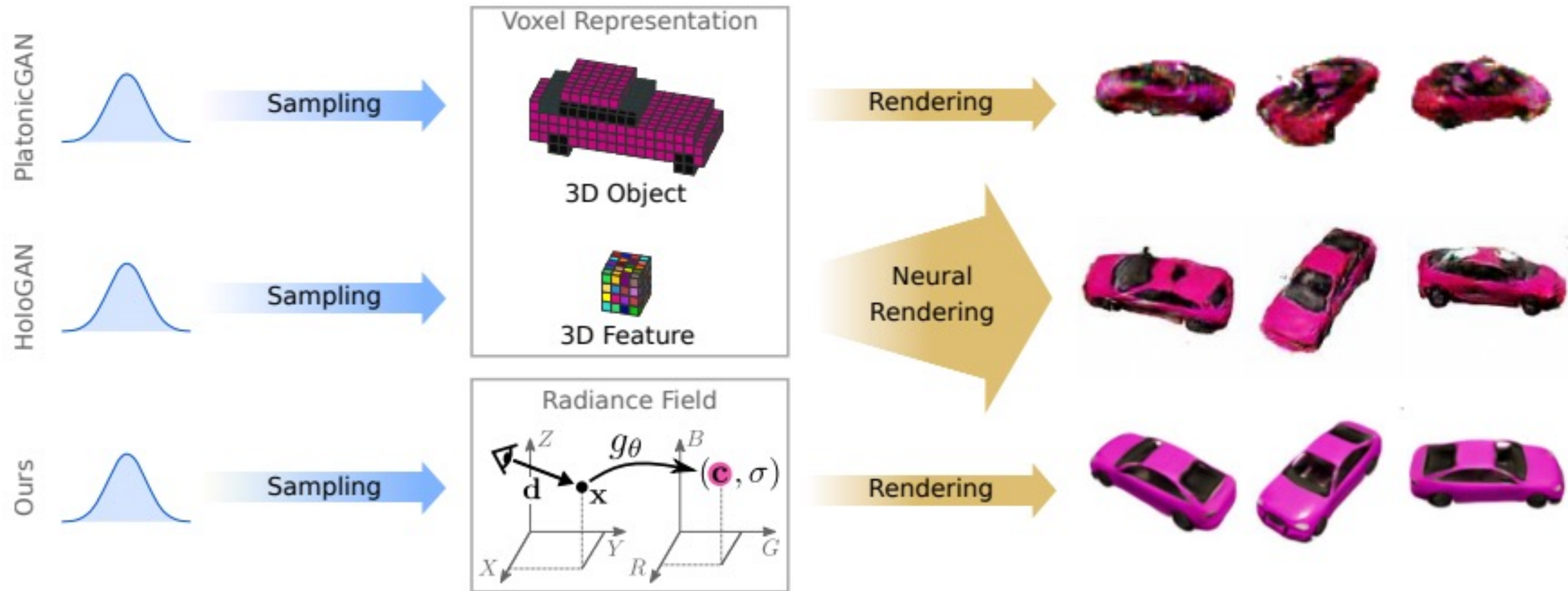
Limitations:

- Do not synthesize geometric outputs (e.g., voxels, SDF).
- No explicit viewpoint consistency. (same issue with Visual Object Networks)

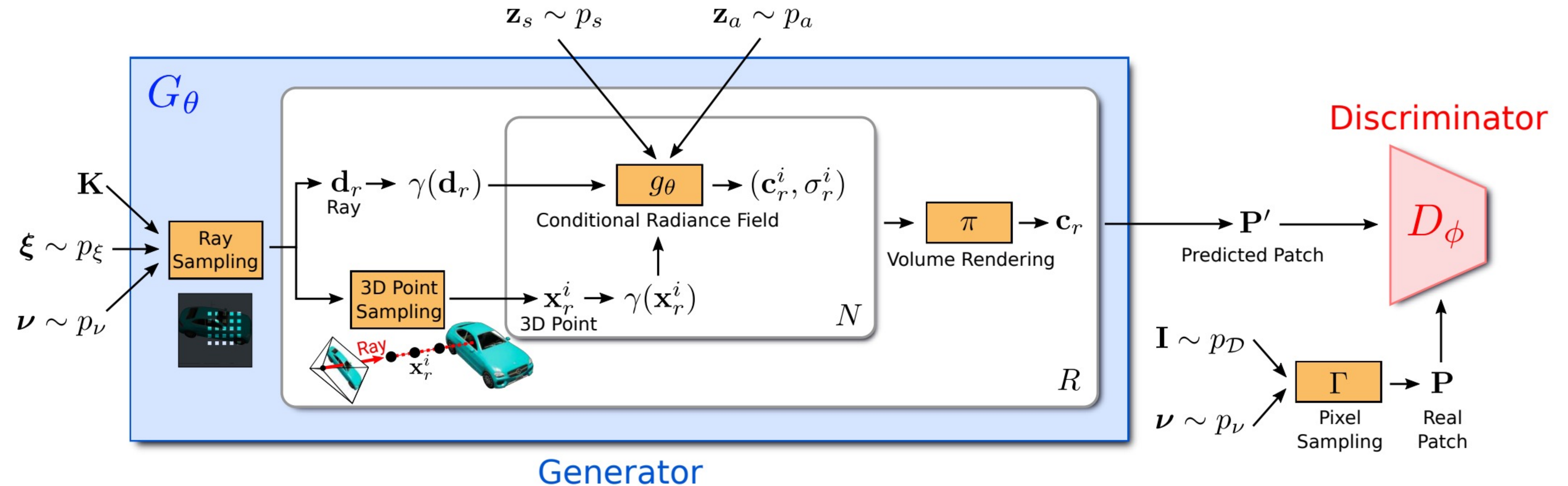
NeRF + GANs

(Neural rendering + Generative Models)

GRAF: Generative Radiance Fields

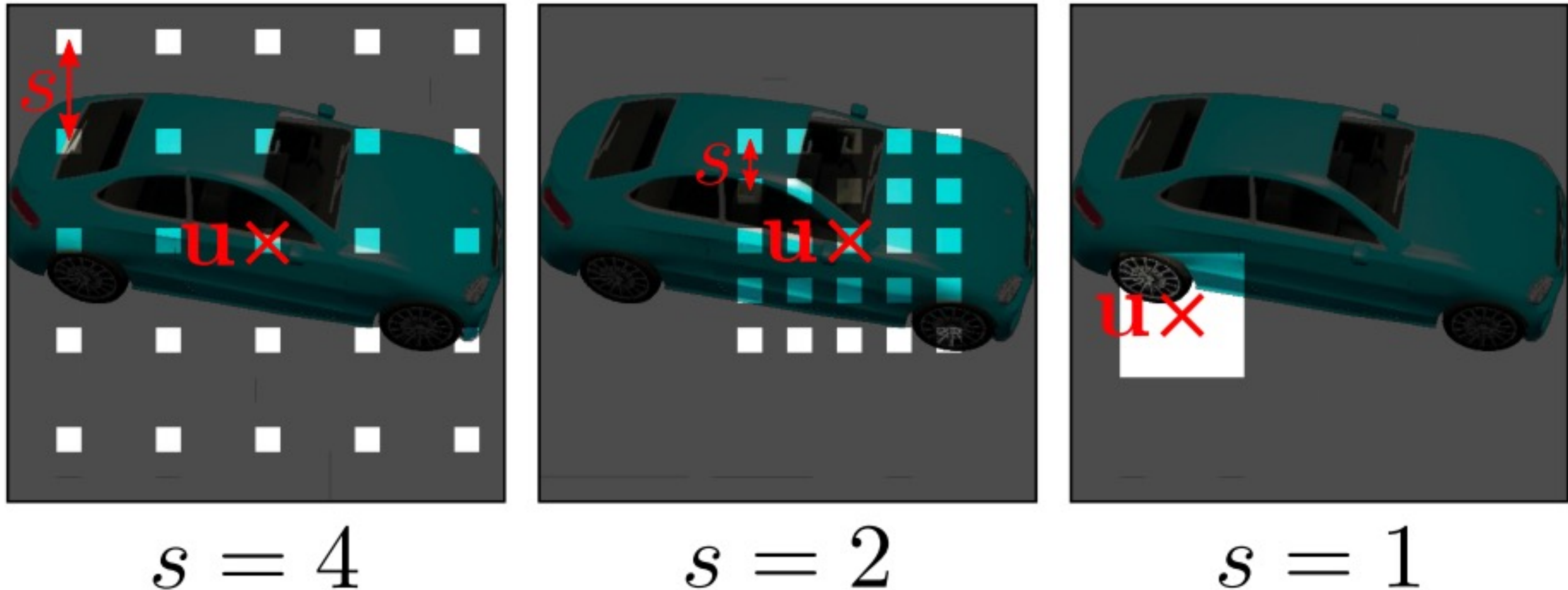


GRAF: Generative Radiance Fields



- **NeRF Generator** is conditioned on both shape and appearance code.
- **Patch-based Discriminator** (full-image discriminator is too slow)

GRAF: Generative Radiance Fields

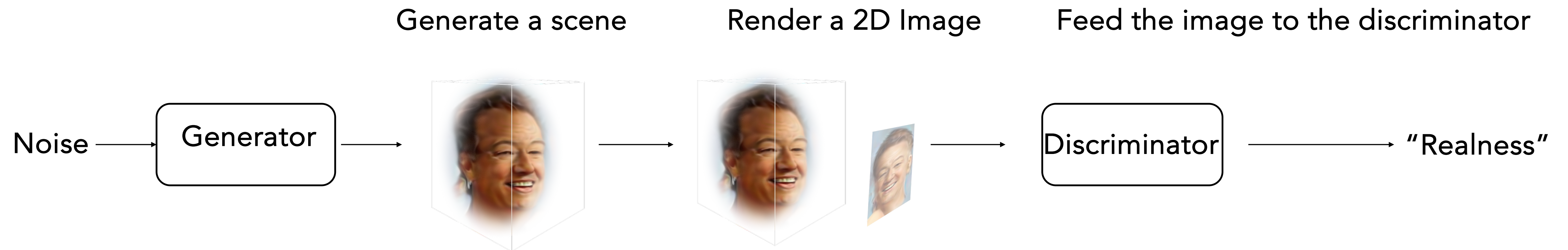


Multi-scale ray sampling

Training a 3D-Aware GAN

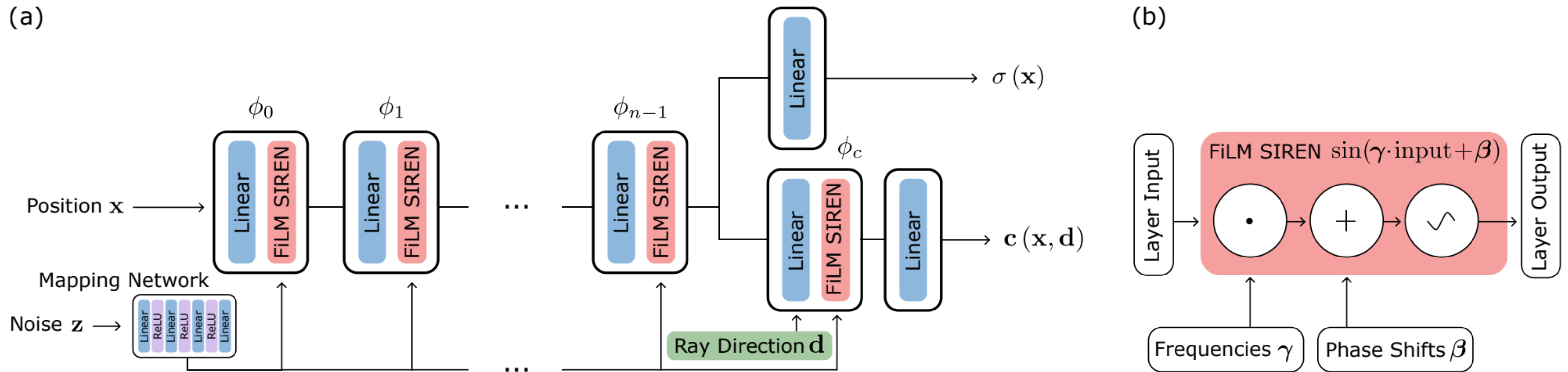
3D-Aware GAN Training Steps

1. Generate a representation of a *scene*
2. Render the scene from a random camera pose
3. Feed the image to a 2D discriminator
4. Backpropagate through the discriminator and differentiable rendering



Slide credit: Eric Chan

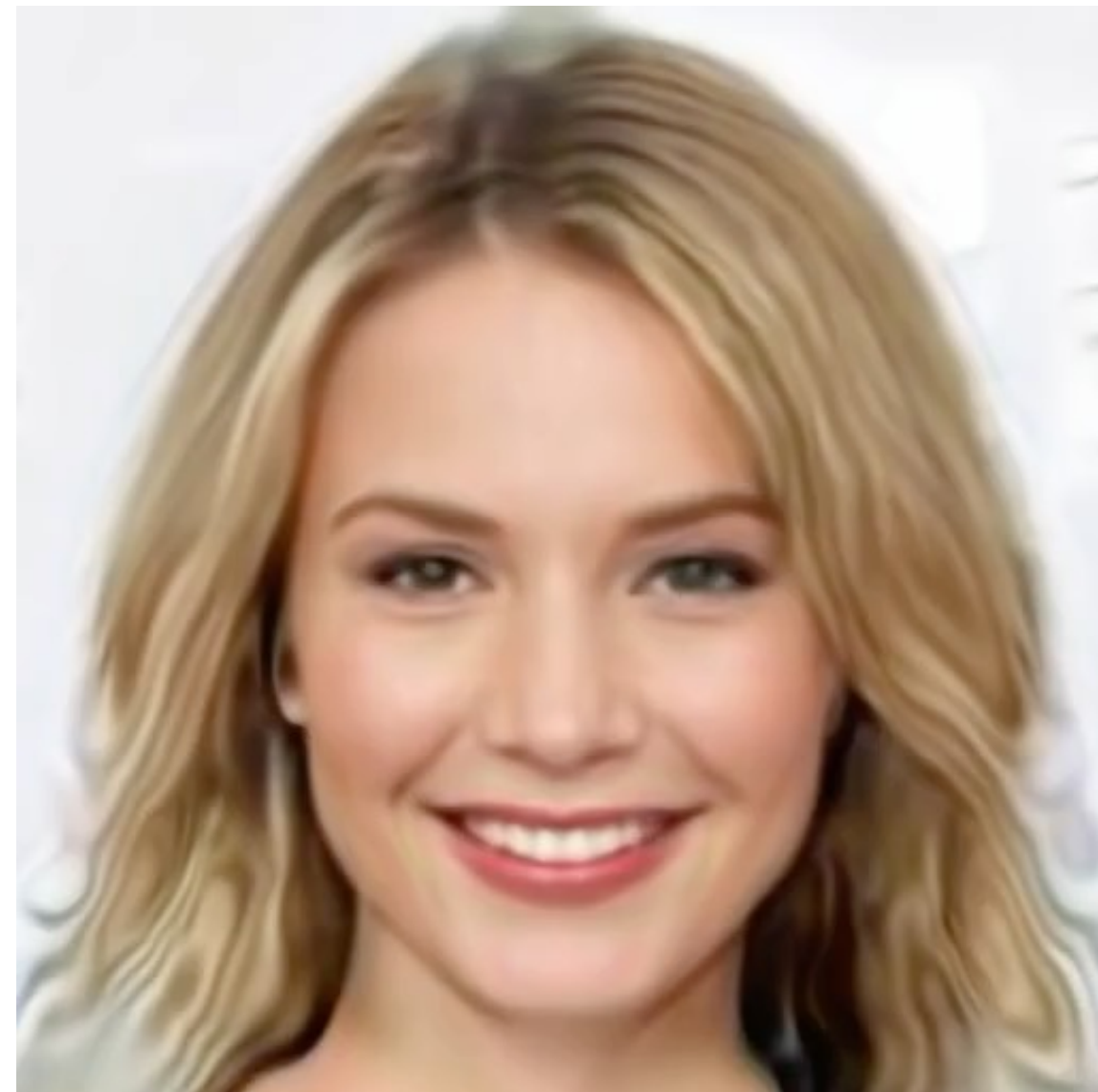
π -GAN



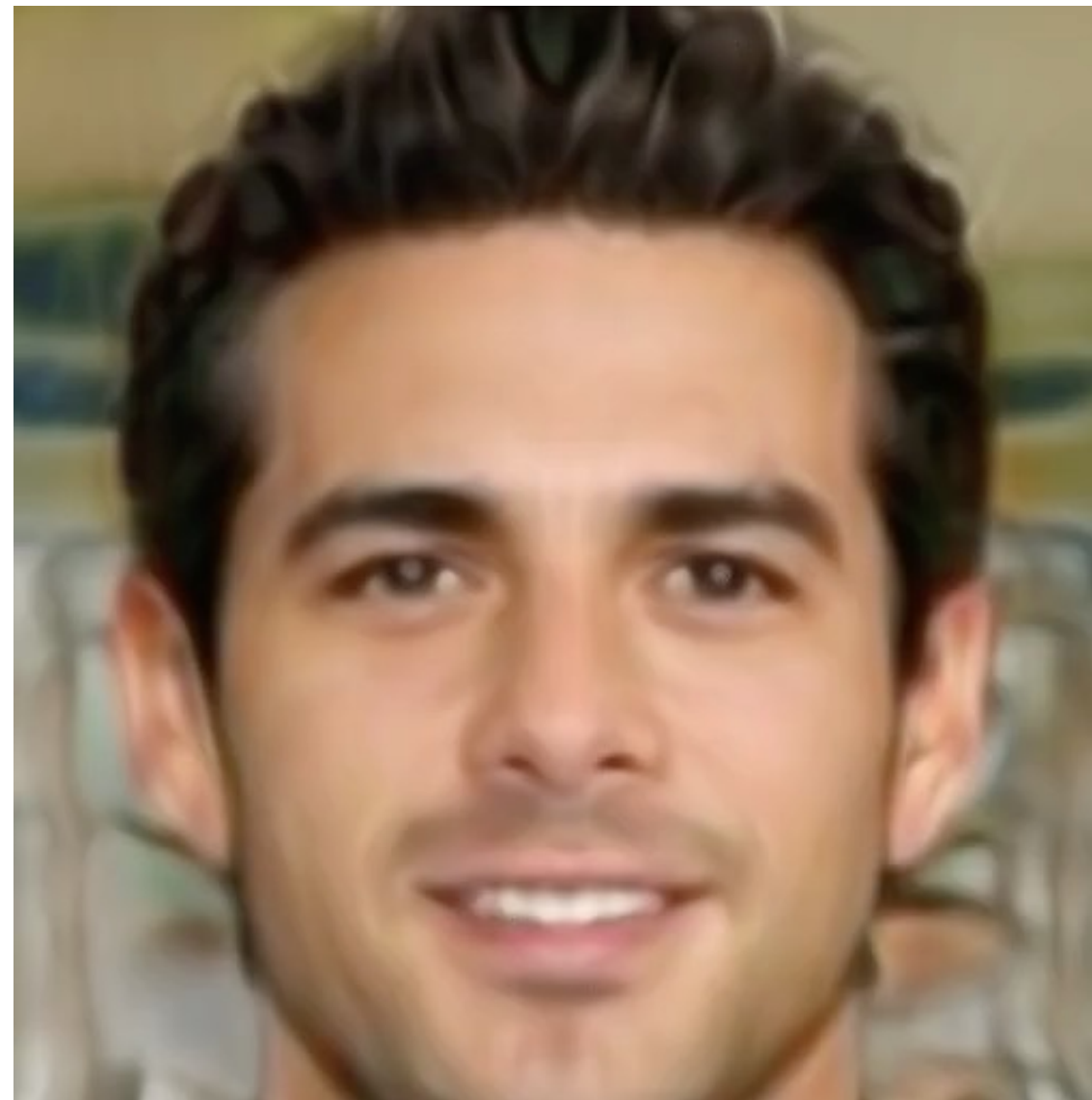
Mapping network + AdaIN (FiLM) + learnable positional encoding

$$\phi_i(\mathbf{x}_i) = \sin(\gamma_i \cdot (\mathbf{W}_i \mathbf{x}_i + \mathbf{b}_i) + \beta_i)$$

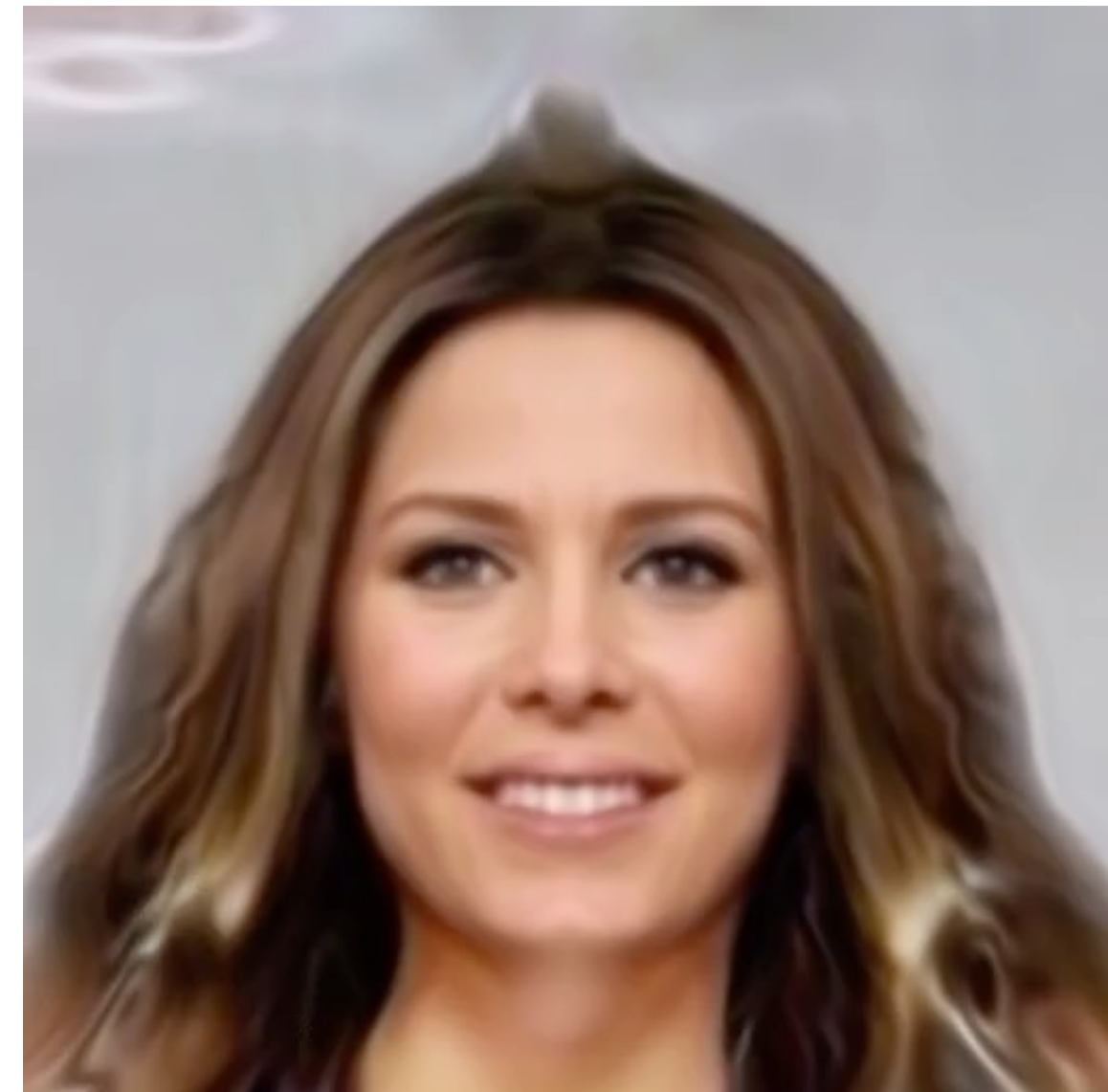
π -GAN



Focal Length



Camera Position



Latent Interpolation

Slide credit: Eric Chan

π -GAN: Periodic Implicit Generative Adversarial Networks for 3D-Aware Image Synthesis [Chan et al., 2021]

π -GAN



Slide credit: Eric Chan

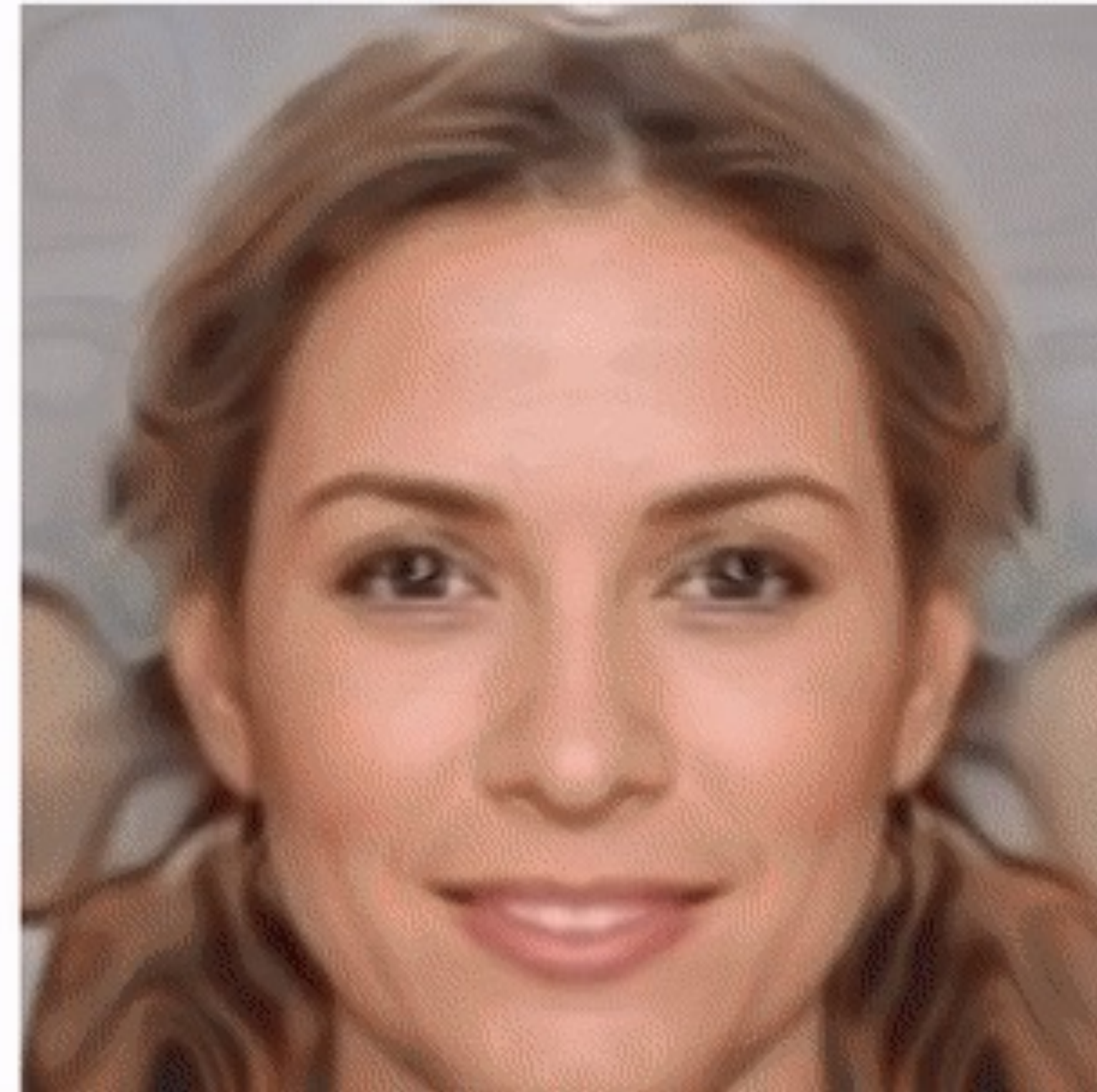
π -GAN: Periodic Implicit Generative Adversarial Networks for 3D-Aware Image Synthesis [Chan et al., 2021]

π -GAN

Target



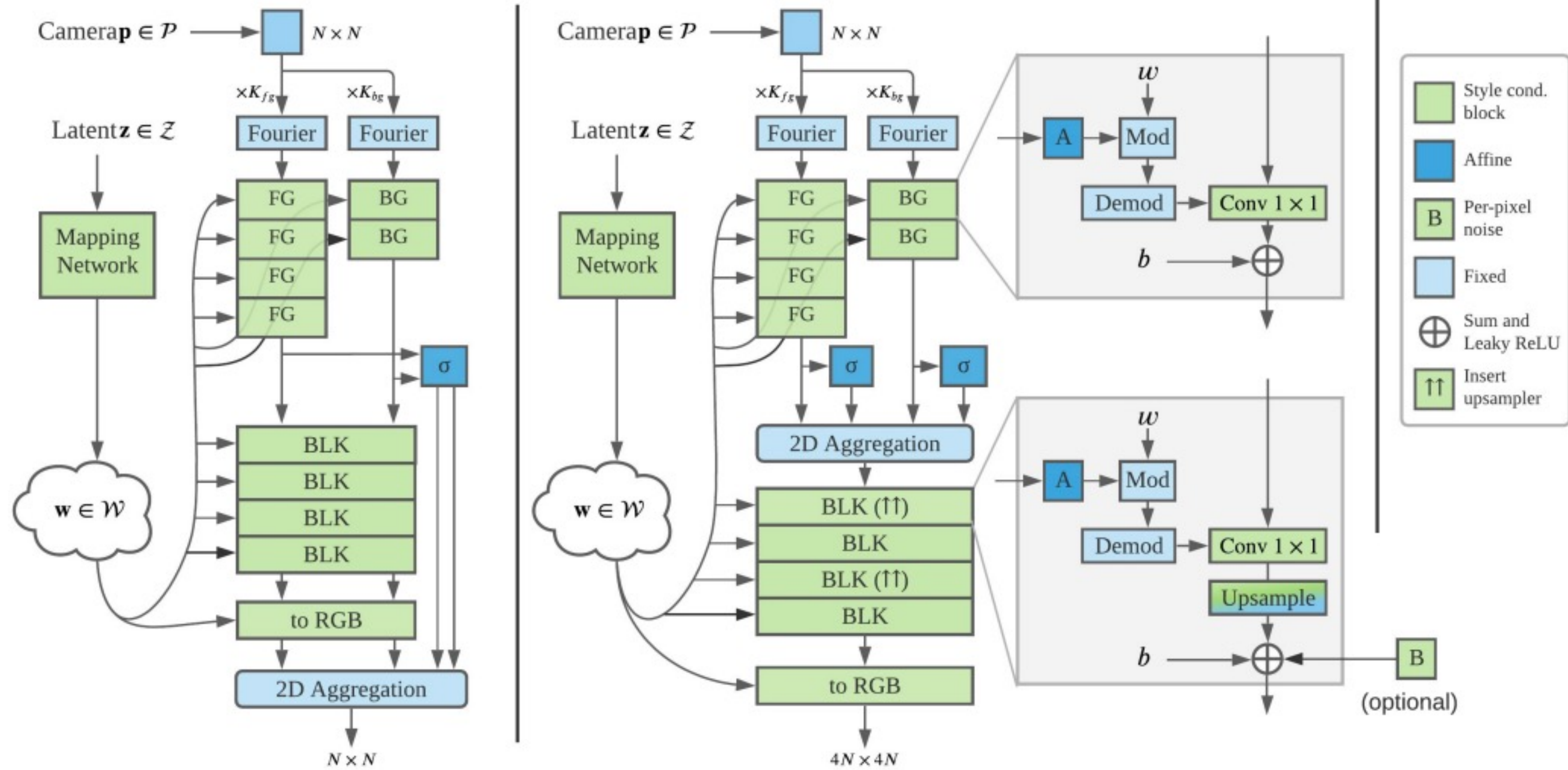
Reconstruction



Slide credit: Eric Chan

π -GAN: Periodic Implicit Generative Adversarial Networks for 3D-Aware Image Synthesis [Chan et al., 2021]

Advanced Architectures: StyleNeRF



Baseline architecture

Proposed architecture

rendering features via volumetric rendering + GANs-based upsampler

Also see recent work: e.g., StyleNeRF [Gu et al.], EG3D [Chan et al.], StyleSDF [Or-EI et al.], ShadeGAN [Pan et al.], ...

StyleNeRF: A Style-based 3D-Aware Generator for High-resolution Image Synthesis [Gu et al., 2021]

MODEL NAME
FFHQ512

CHECKPOINT PATH

TRUNCATION TRICK
0.7

SEED1
4

SEED2
9

LINEAR MIXING RATIO (GEOMETRY)
0

LINEAR MIXING RATIO (APPARENCE)
0


YAW
0

PITCH
0

FOV
12

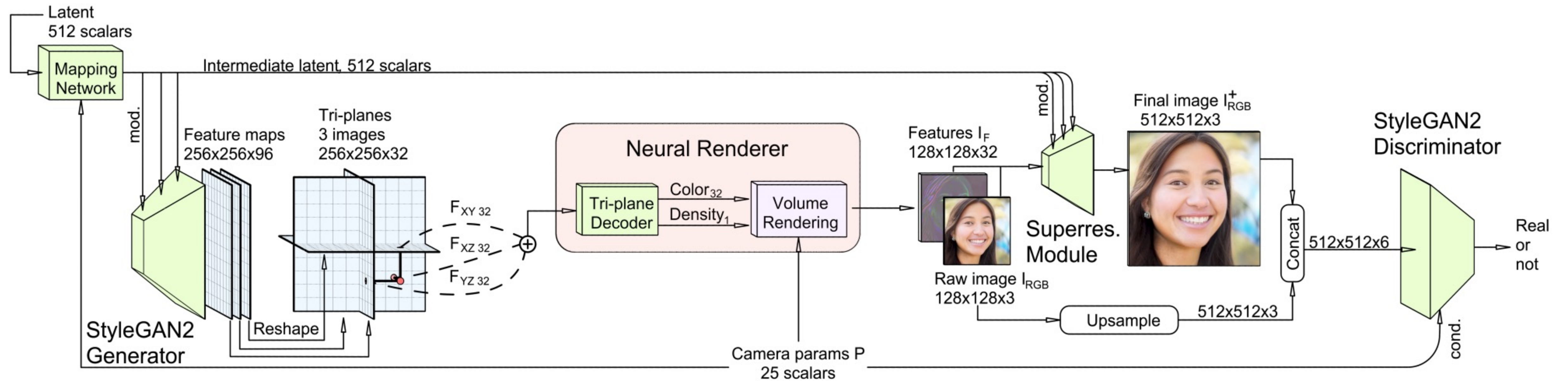
Clear

OUTPUT 0.12s



Screenshot Flag

Advanced Architectures: EG3D (StyleNeRF+Triplane)



**Tri-plane representation
for speed-up**

$$F(x, y, x) \rightarrow F(x, y) + F(x, z) + F(y, z)$$

**Rendering features
via volumetric rendering**

features are useful for upsampling

**Generate final output
via image encoder**

If the model is too slow,
use GAN-based upsampler

Also see recent work: e.g., StyleNeRF [Gu et al.], EG3D [Chan et al.], StyleSDF [Or-EI et al.], ShadeGAN [Pan et al.], ...

EG3D: Efficient Geometry-aware 3D Generative Adversarial Networks [Chan et al., 2021]



EG3D: Efficient Geometry-aware 3D Generative Adversarial Networks [Chan et al., 2021]

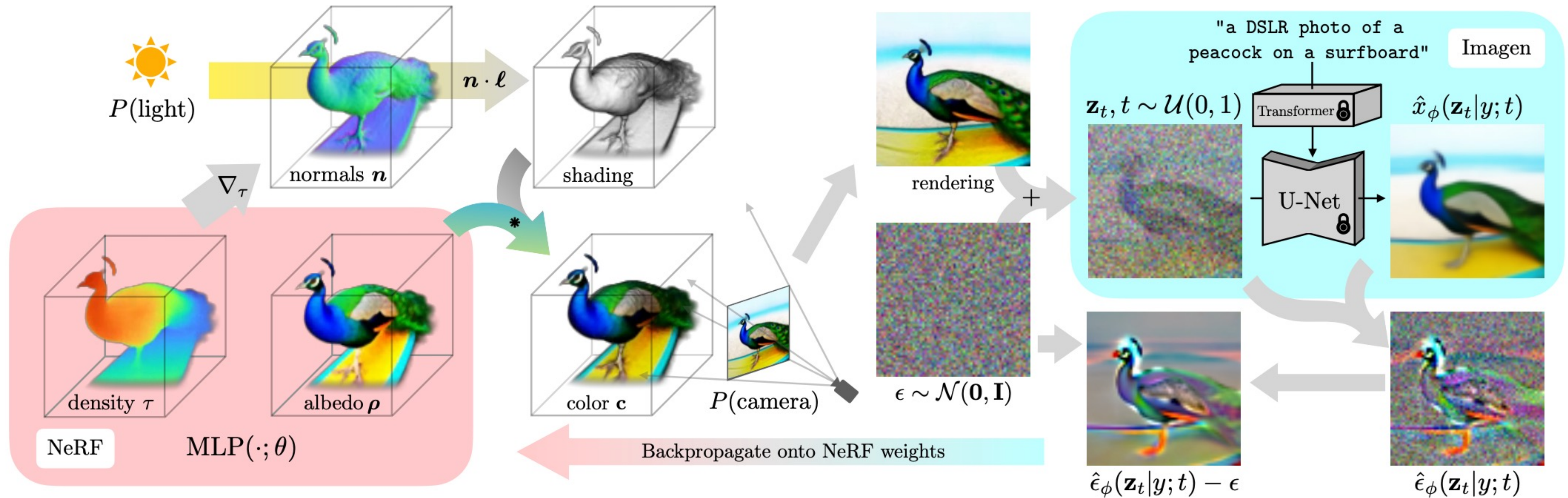
Text-based Editing with Generative NeRFs

Text-based Editing

a DSLR photo of a squirrel
wearing a purple hoodie
reading a book



Text-based Editing



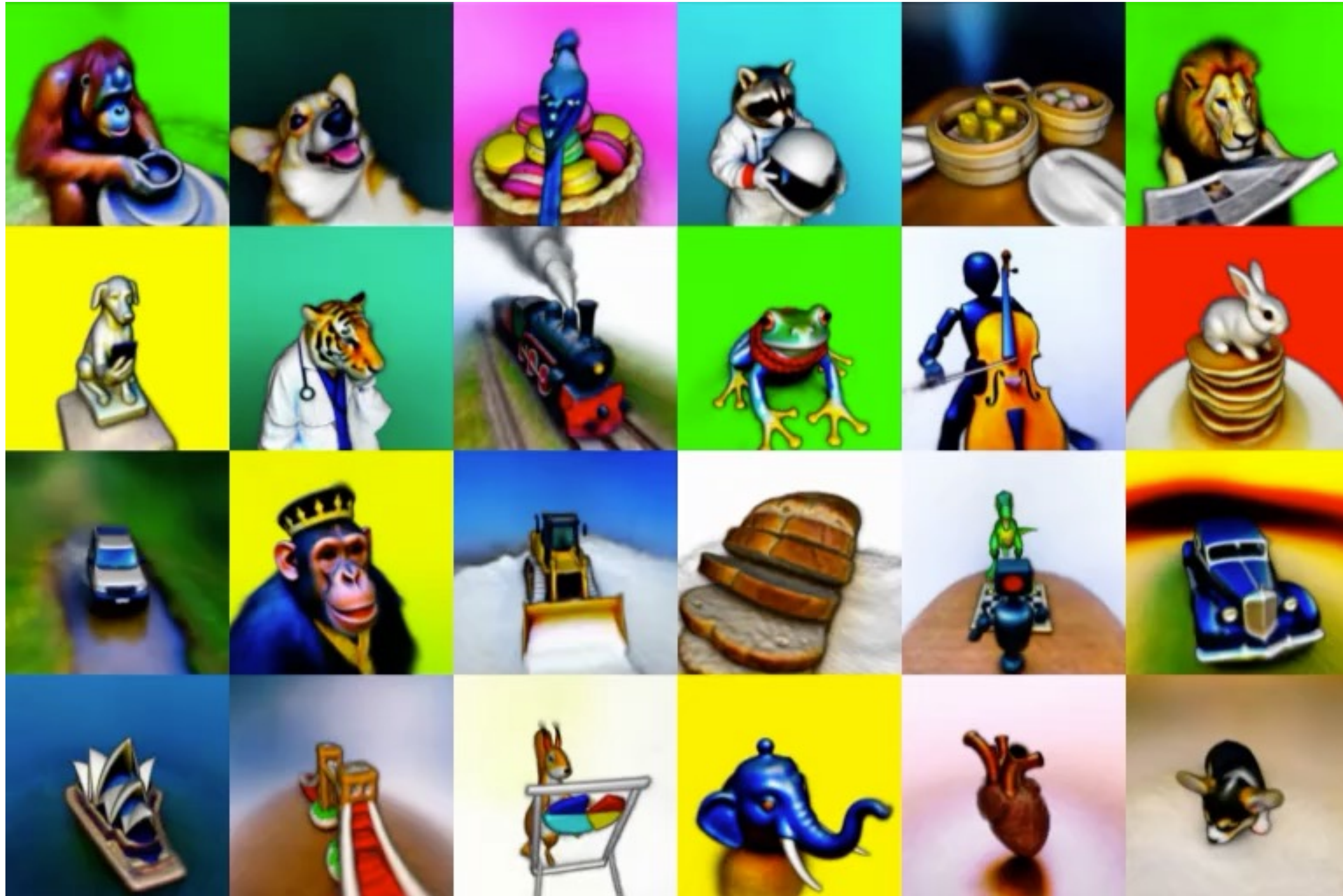
FOR loop

Step 1. Render a view using existing NeRF

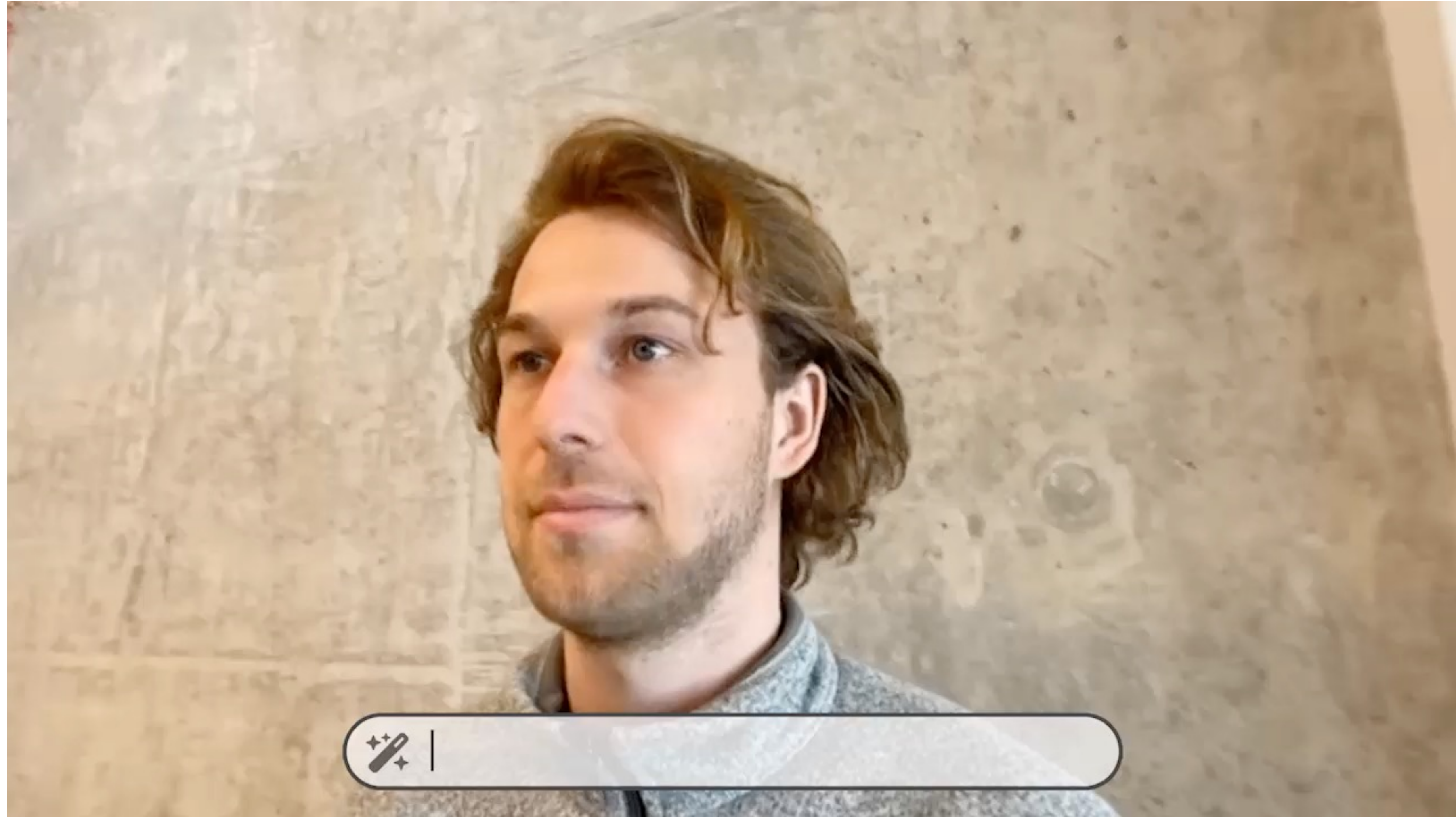
Step 2. Add noise and denoise using a pre-trained Stable Diffusion model

Step 3. Update NeRF parameters with the gradient (difference between added and predicted noises)

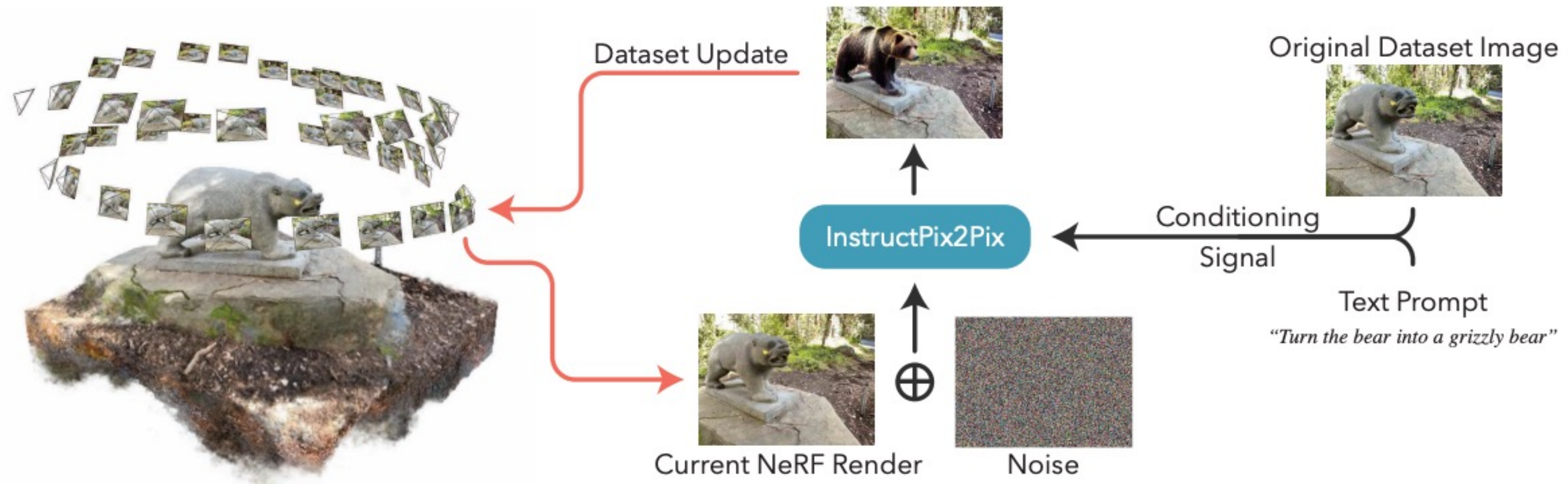
Text-based Editing



Instruct NeRF2NeRF



Instruct NeRF2NeRF



FOR loop

Step 1. Render a view using existing NeRF

Step 2. Use InstructPix2Pix to produce output images

Step 3. Update NeRF parameters with the generated result from Step 2

InstructPix2pix: image-conditional diffusion model (<https://www.timothybrooks.com/instruct-pix2pix/>)

Thank You!

<https://learning-image-synthesis.github.io/>