



3D-aware Synthesis (part II)

Jun-Yan Zhu

16-726, Spring 2022

HW3 Class Choice Award Winner

Yutian Lei

Honorable mentions

Violet Han, Tomas Cabezon Pedroso



Virtual Game night (Friday 9 pm)

battle

0:18 team 1

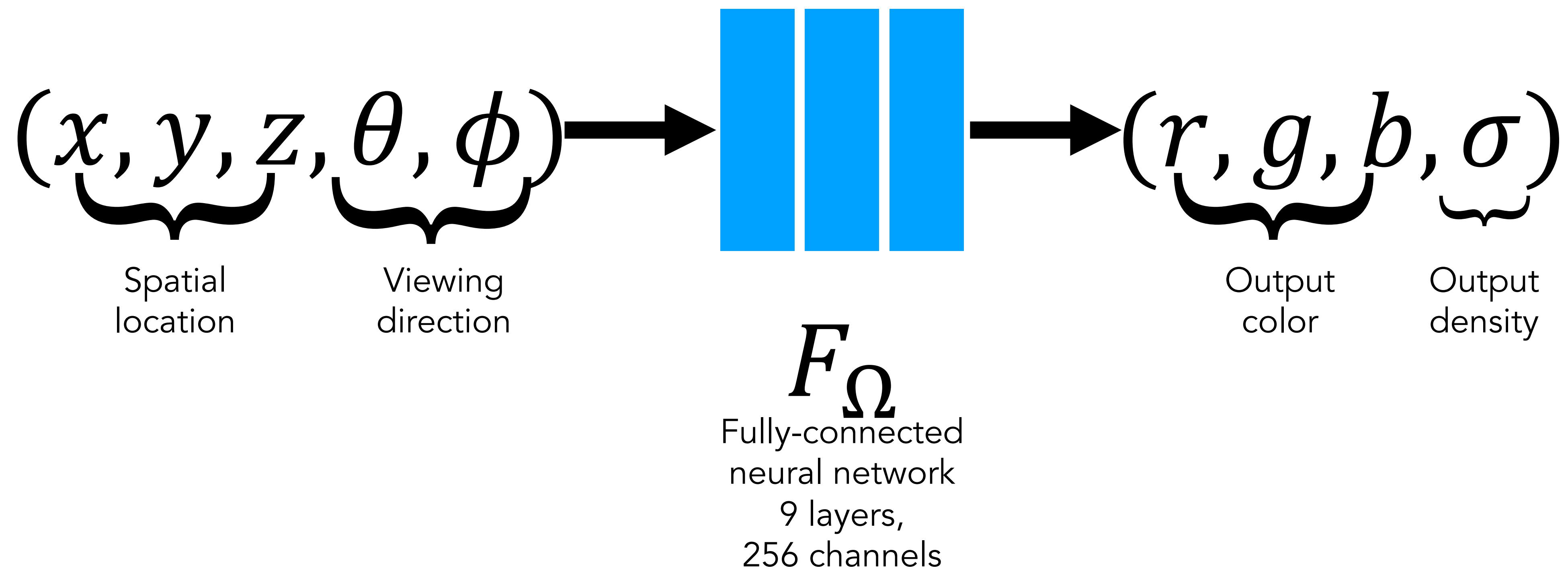
The image shows a digital drawing interface. On the left, there is a white square with a green border containing a black outline drawing of two crossed swords. To the right of the drawing, there is a list of status updates:

- sean is drawing
- ! rita and rajeev are guessing
- ✗ rita guessed sword
- ✗ rajeev guessed duel
- ✓ rita guessed battle

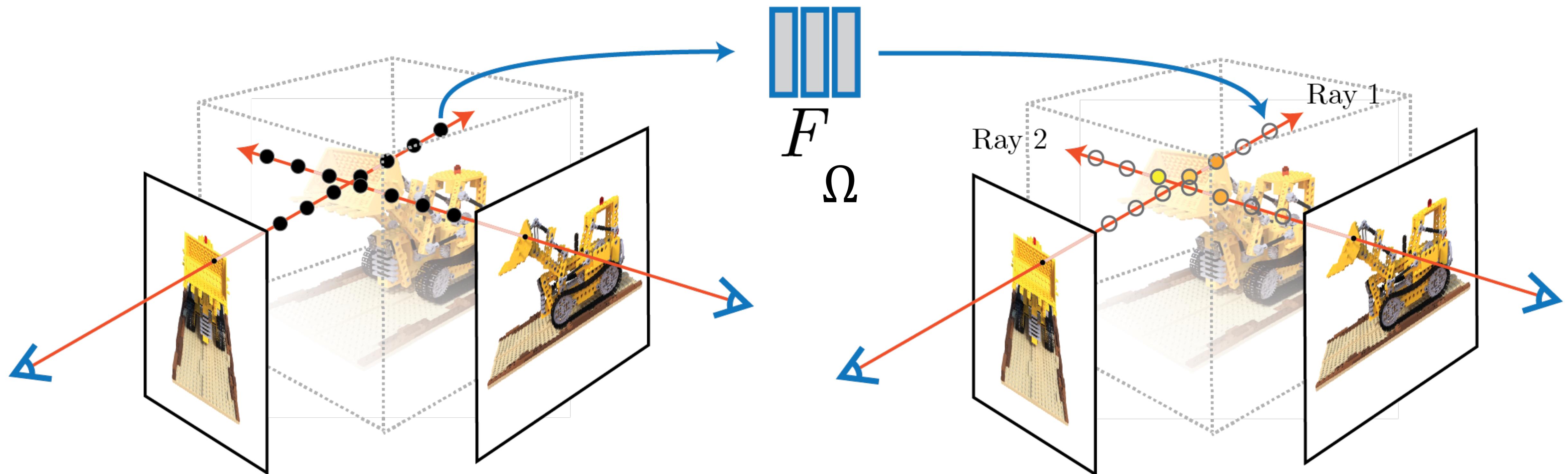
At the bottom, there is a text input field with the placeholder "type your guess..." and a blue button labeled "guess".

NeRF (neural radiance fields):
Neural networks as a volume representation,
using volume rendering to do view
synthesis.(x, y, z, θ, ϕ) \rightarrow *color, opacity*

Representing a scene as a continuous 5D function

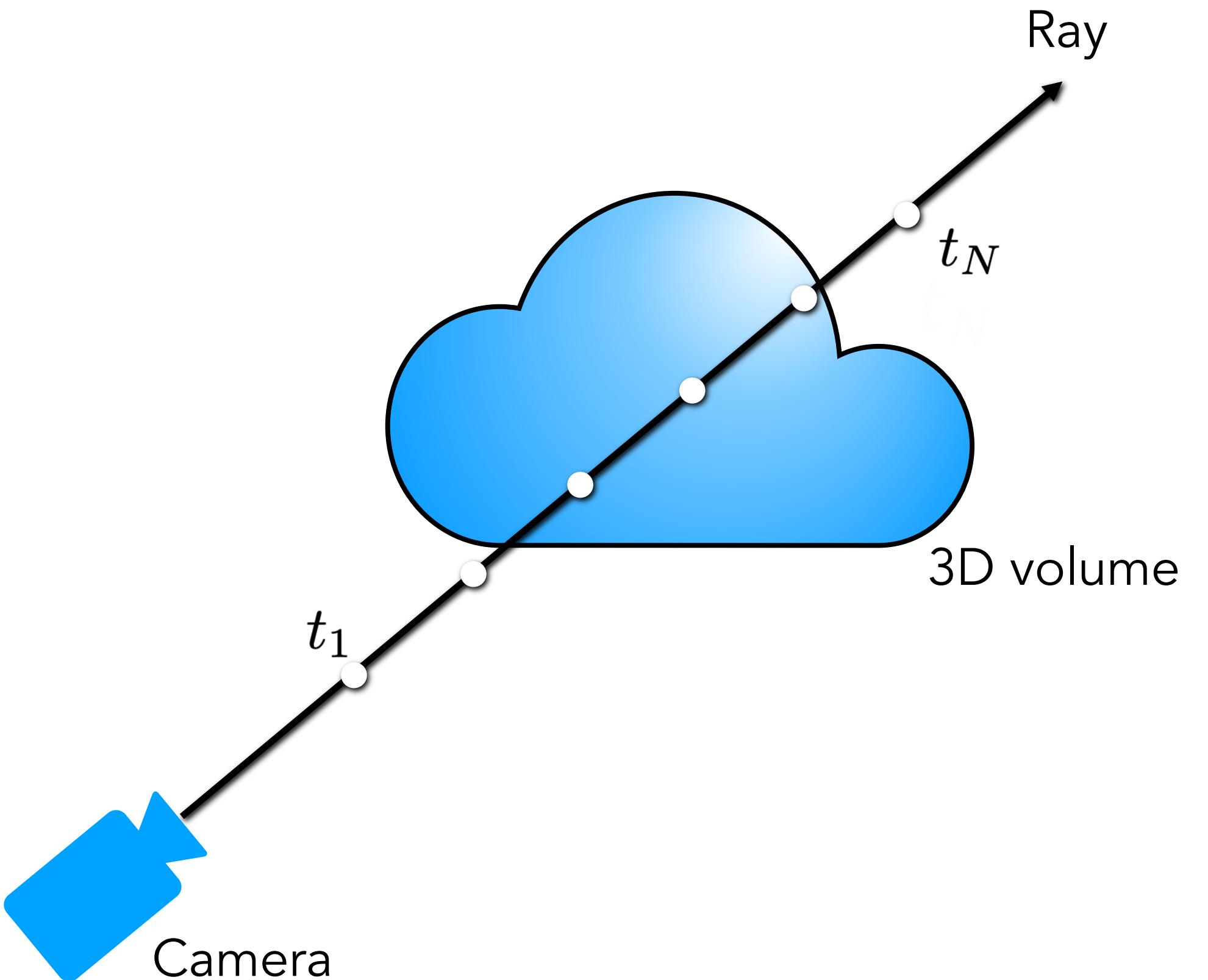


Generate views with traditional volume rendering



Generate views with traditional volume rendering

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

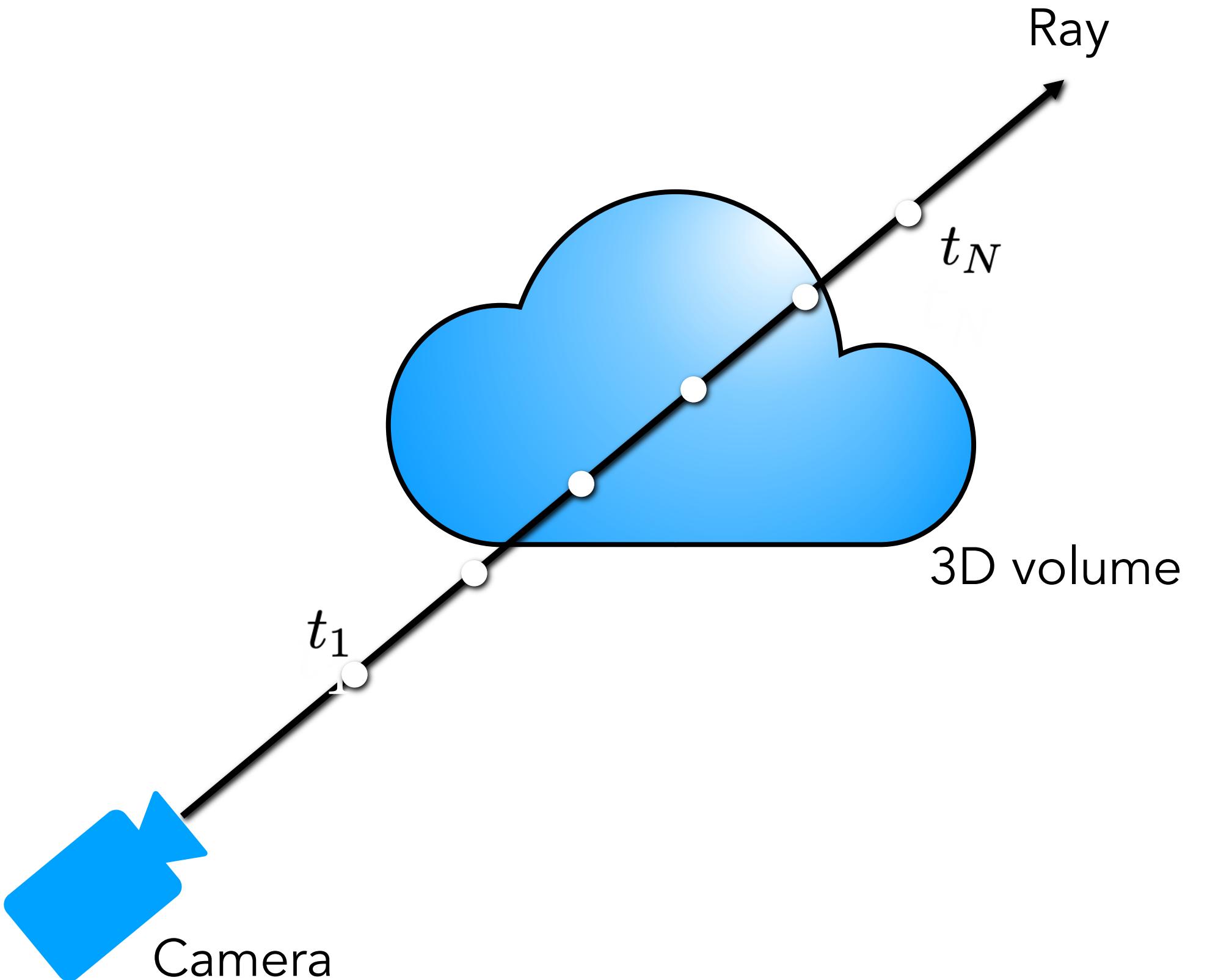


Generate views with traditional volume rendering

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

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

weights colors

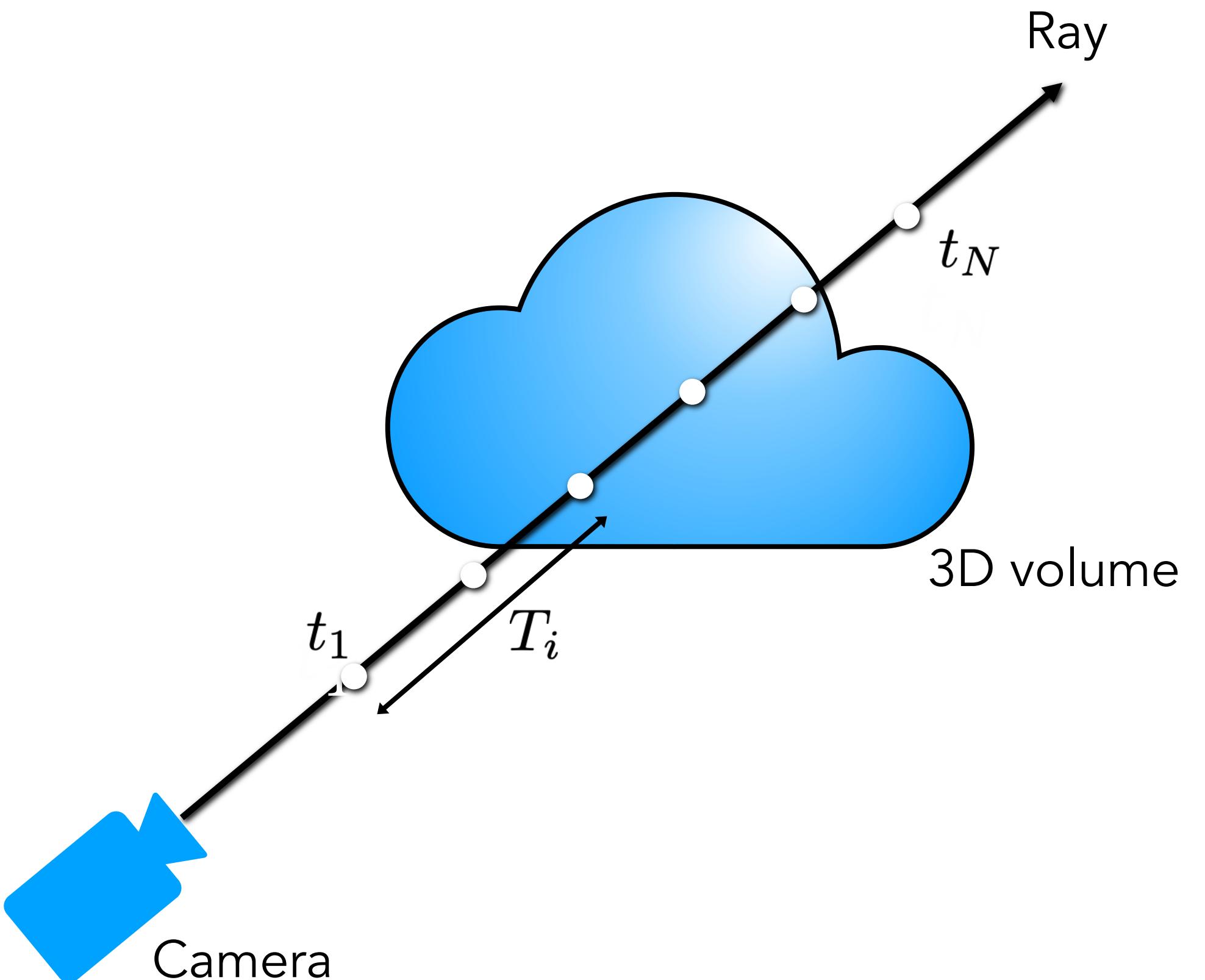


Generate views with traditional volume rendering

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

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 $r(t) = o + td$:

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

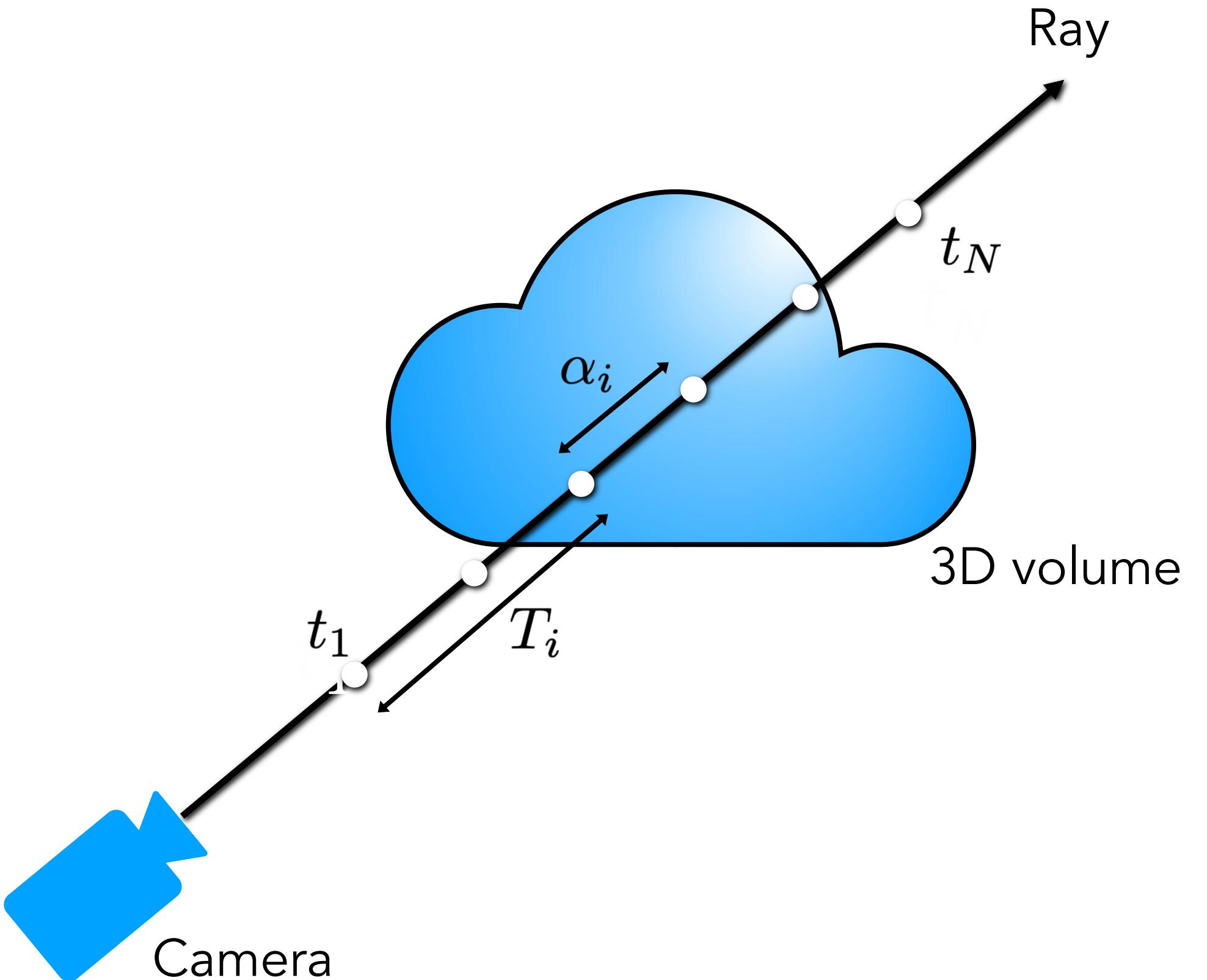
weights colors

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 $r(t) = o + td$:

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

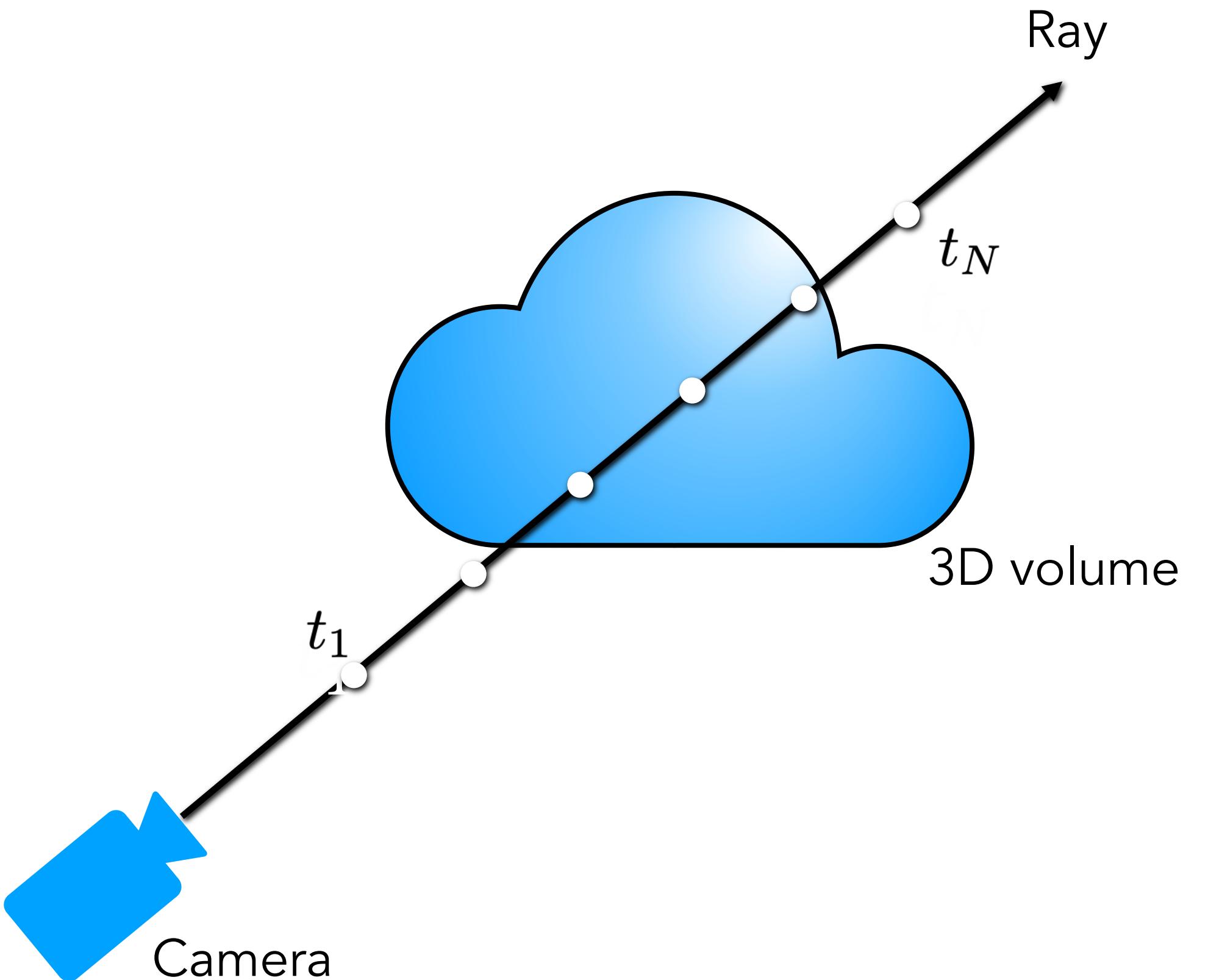
weights colors

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 :

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



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

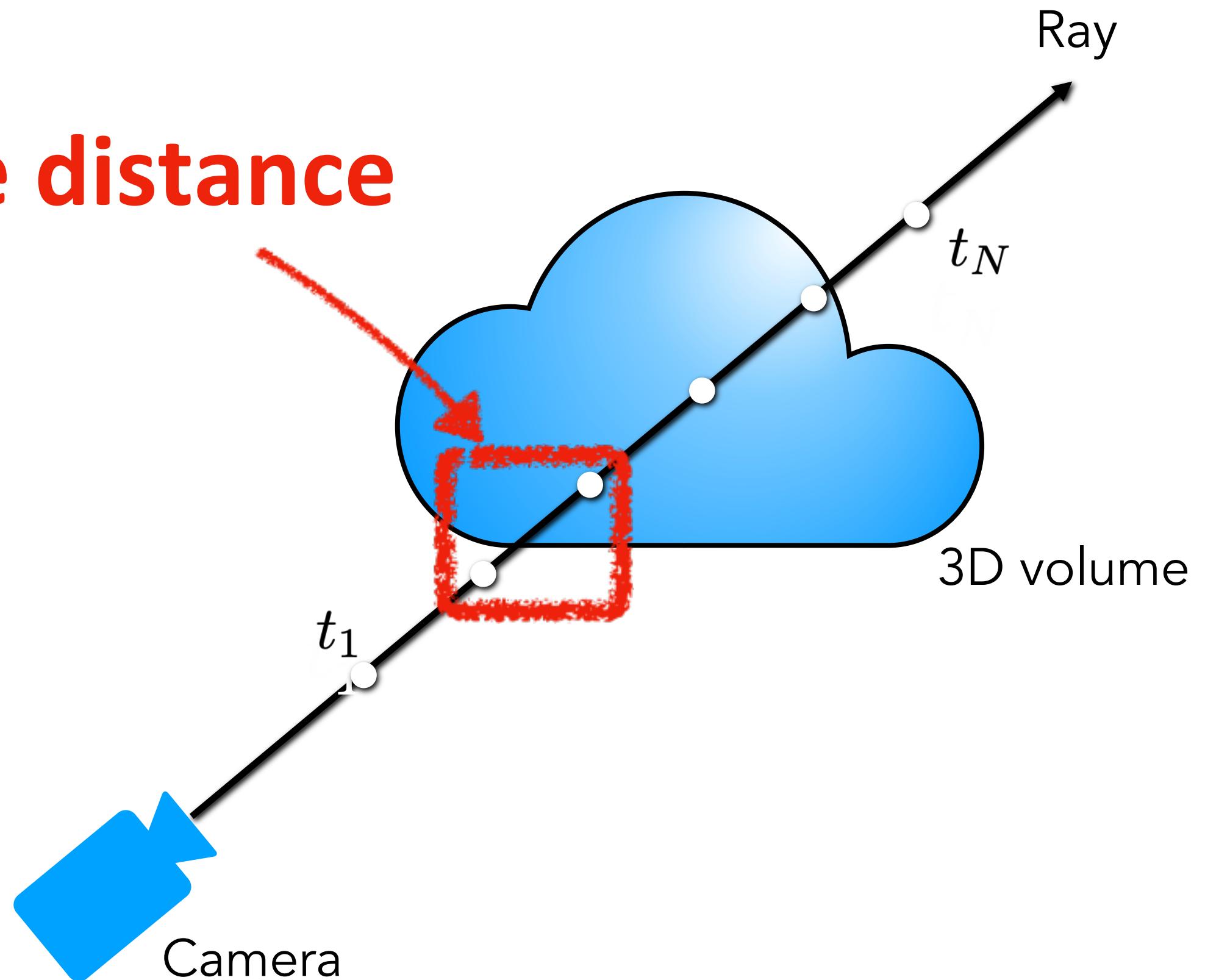
How much light is blocked earlier along ray:

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

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 $r(t) = o + td$:

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

weights

colors

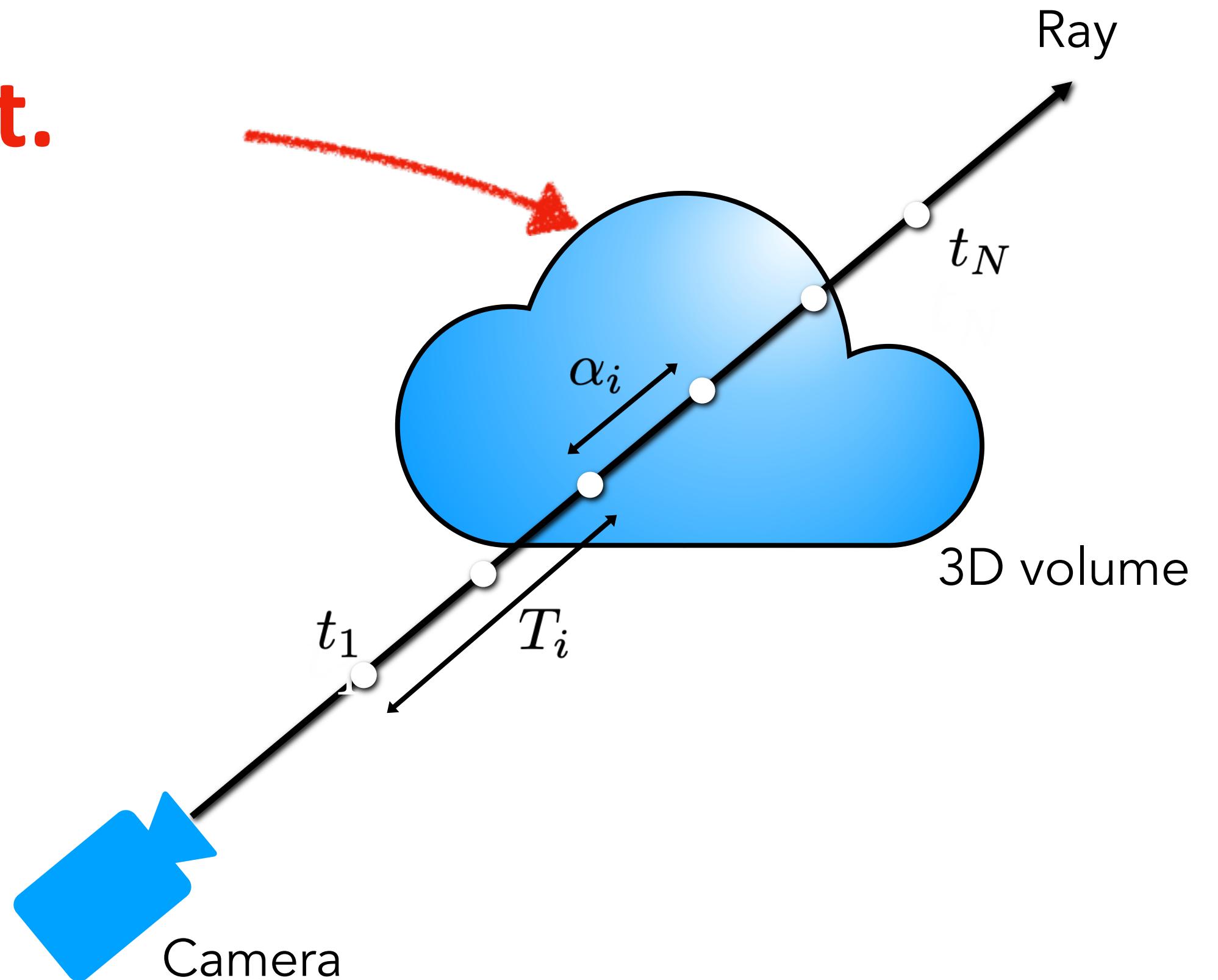
differentiable w.r.t.

How much light is blocked earlier along ray:

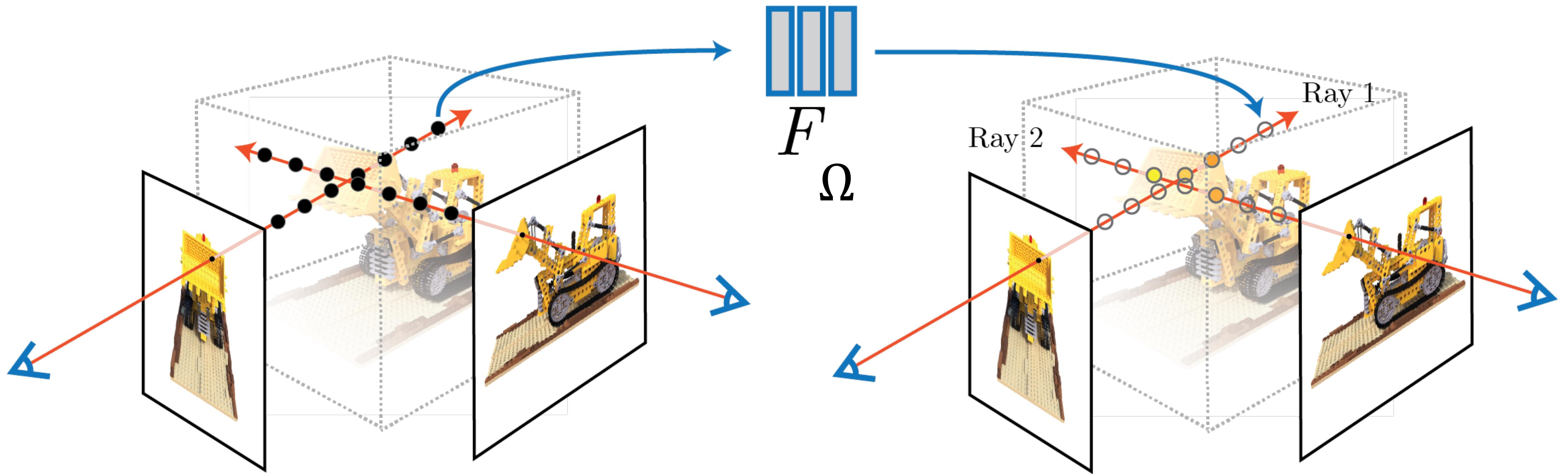
$$T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$$

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

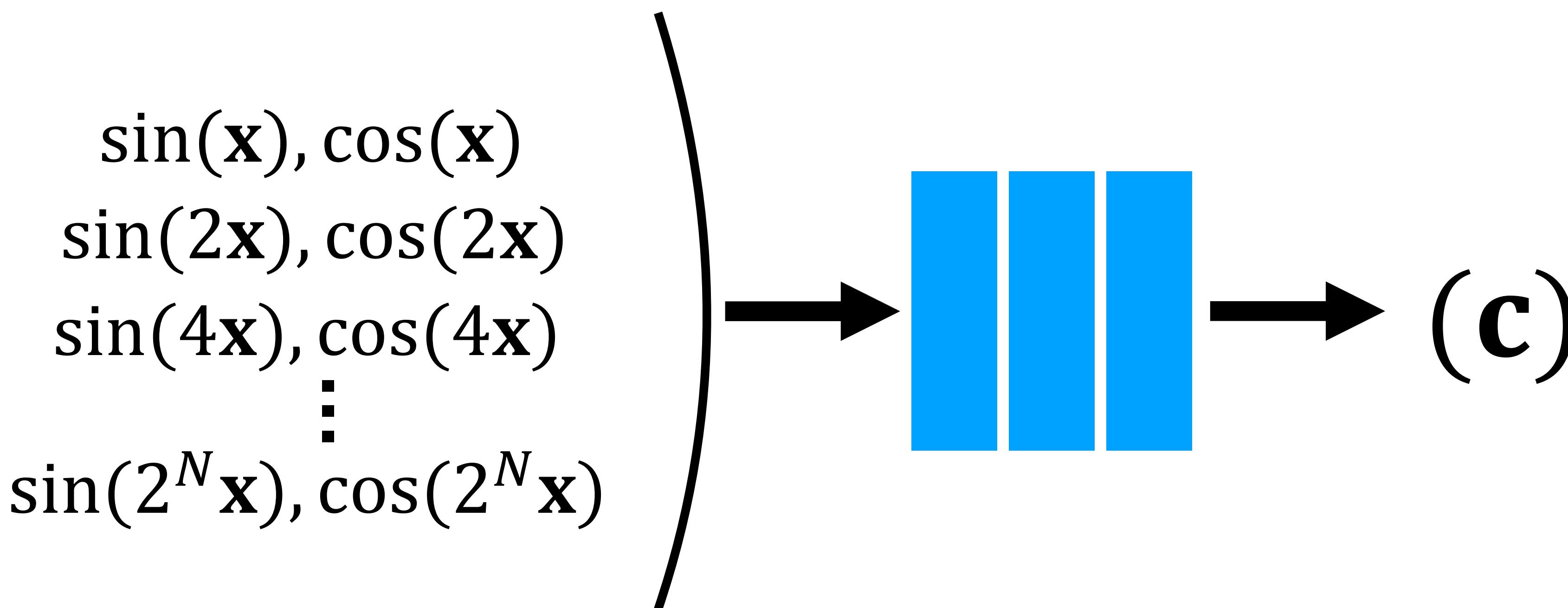


$$\min_{\Omega} \sum_i \| \text{render}^{(i)}(F_\Omega) - I_{\text{gt}}^{(i)} \|^2$$

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

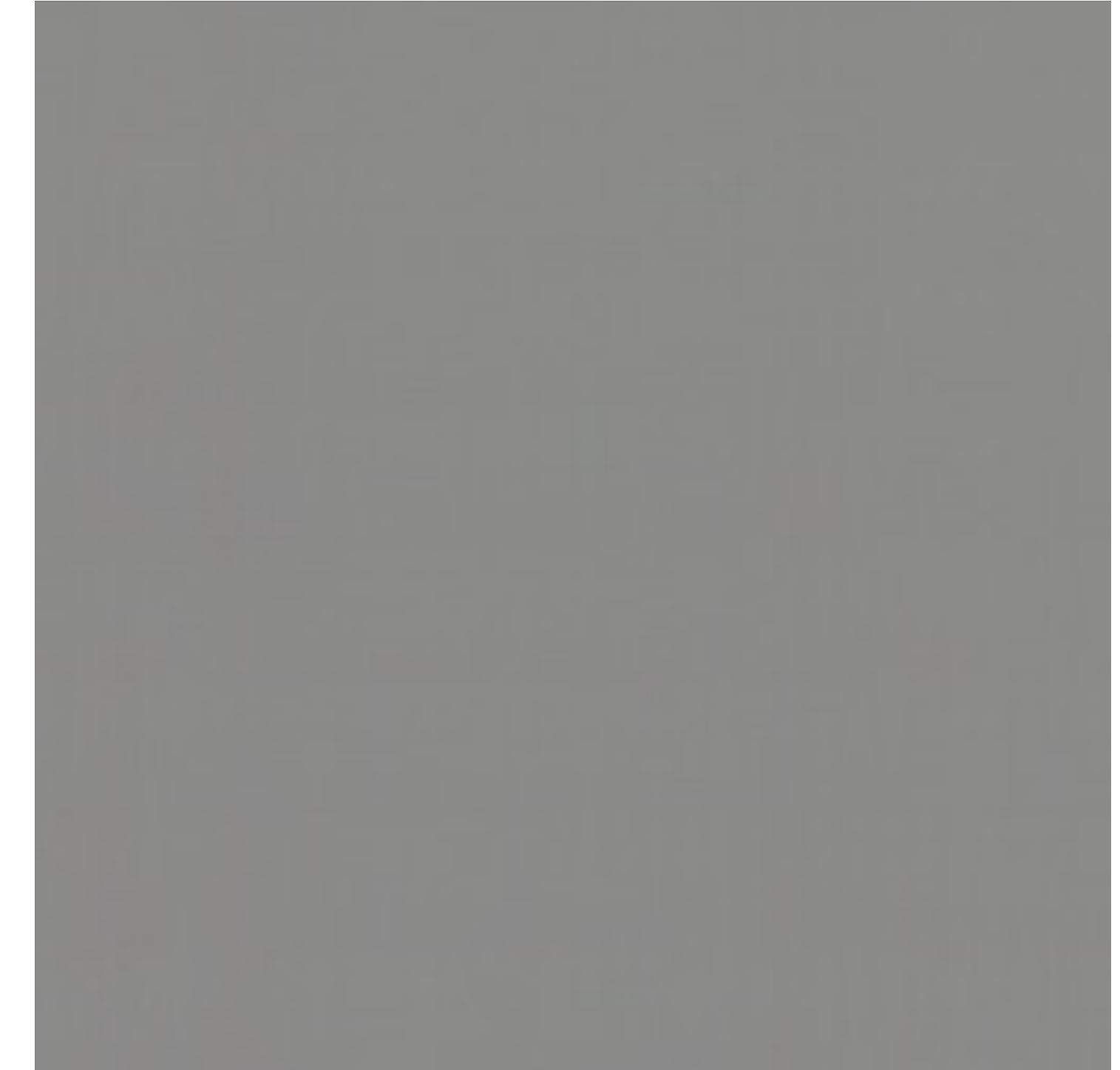
Ground truth image



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

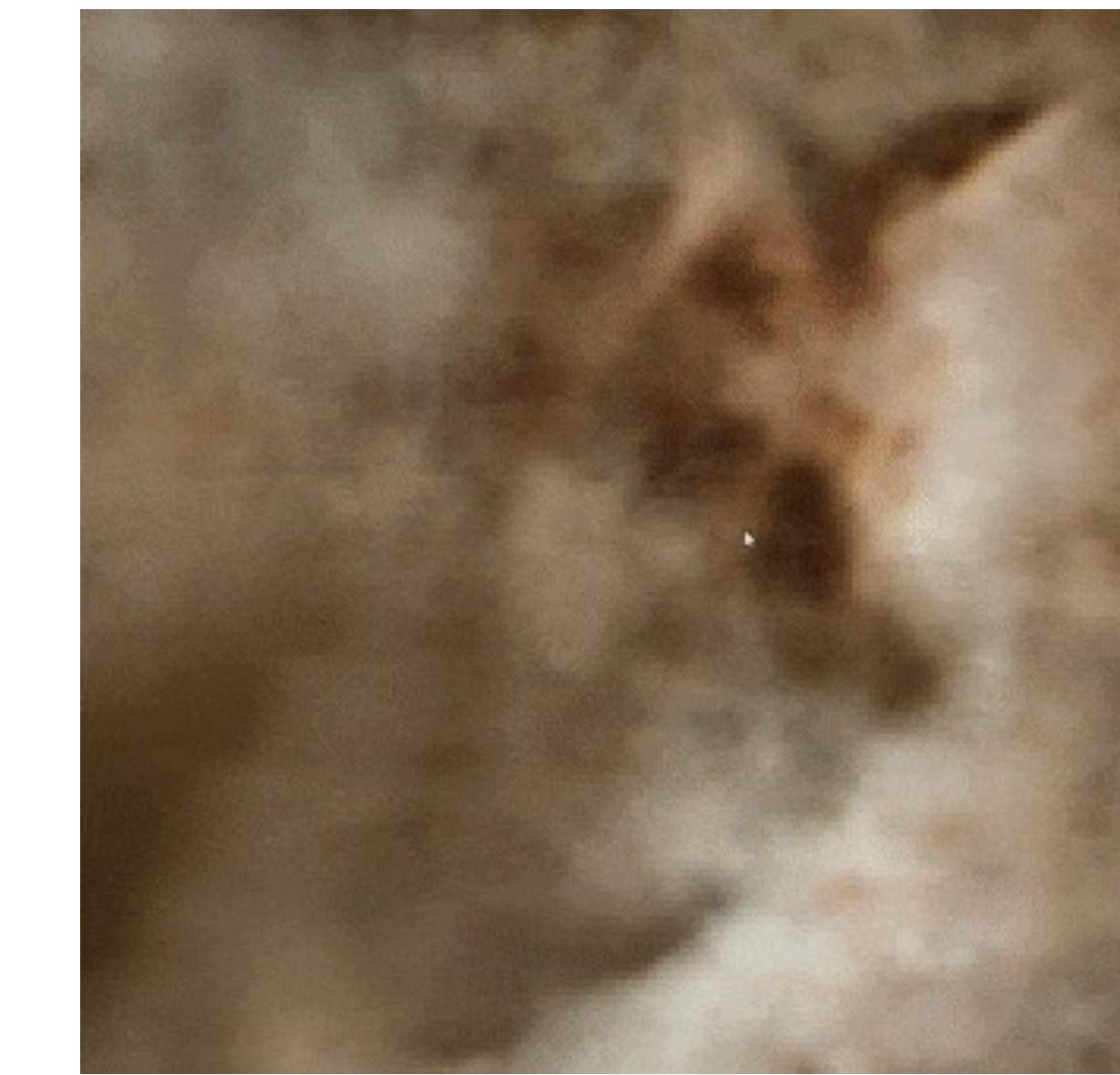
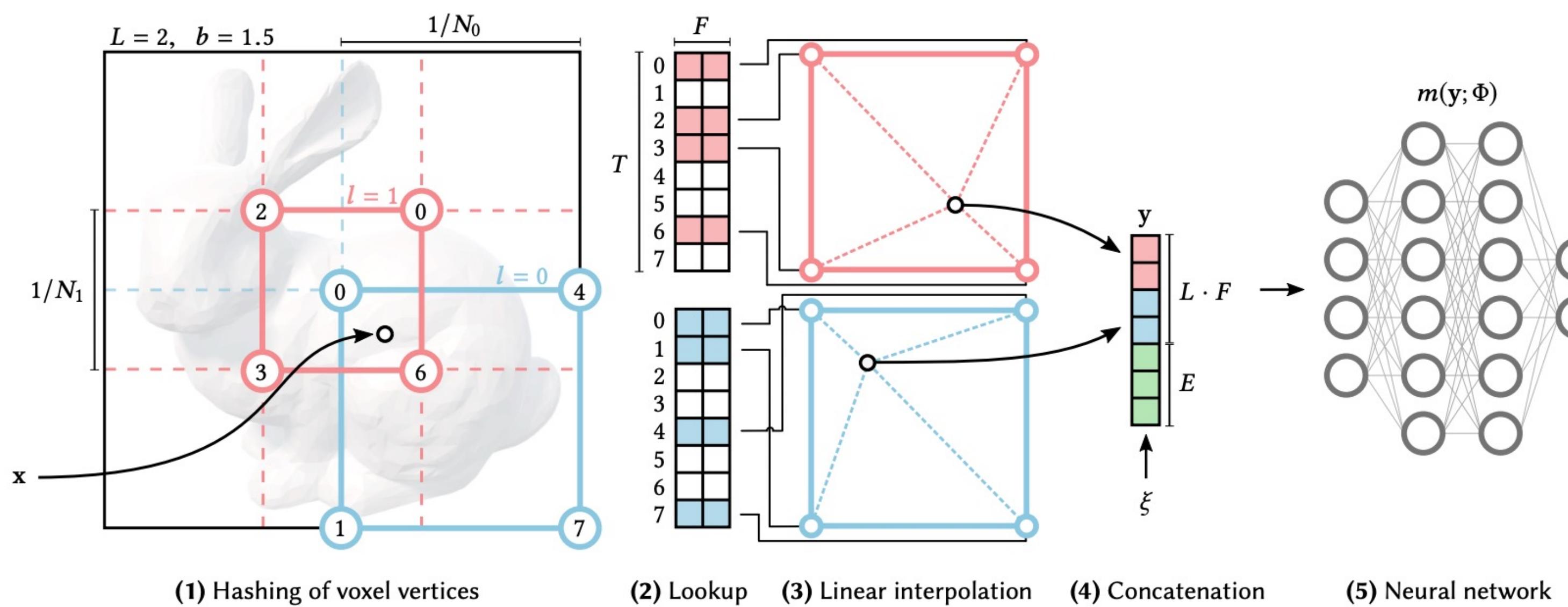
- Use Structure from Motion (e.g., [COLMAP](#)) to initialize camera poses
- Incorrect camera poses lead to bad results
- Joint optimization of camera poses and scene presentation.



Implementation Details

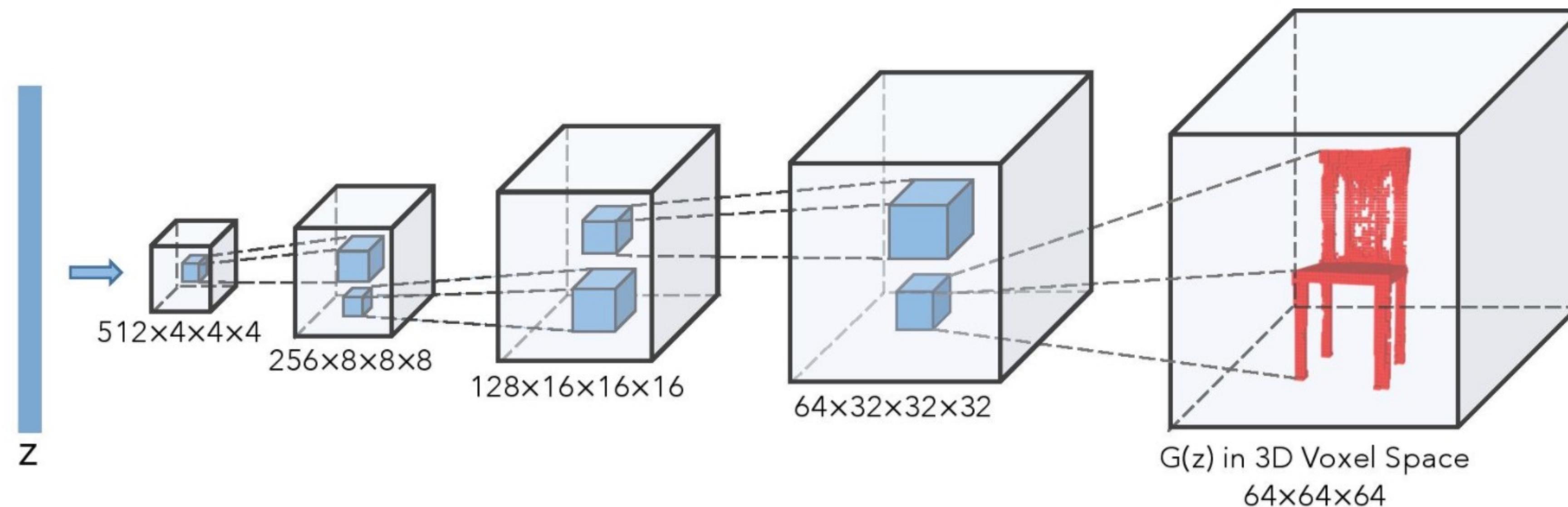
Training and inference speed:

- 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)

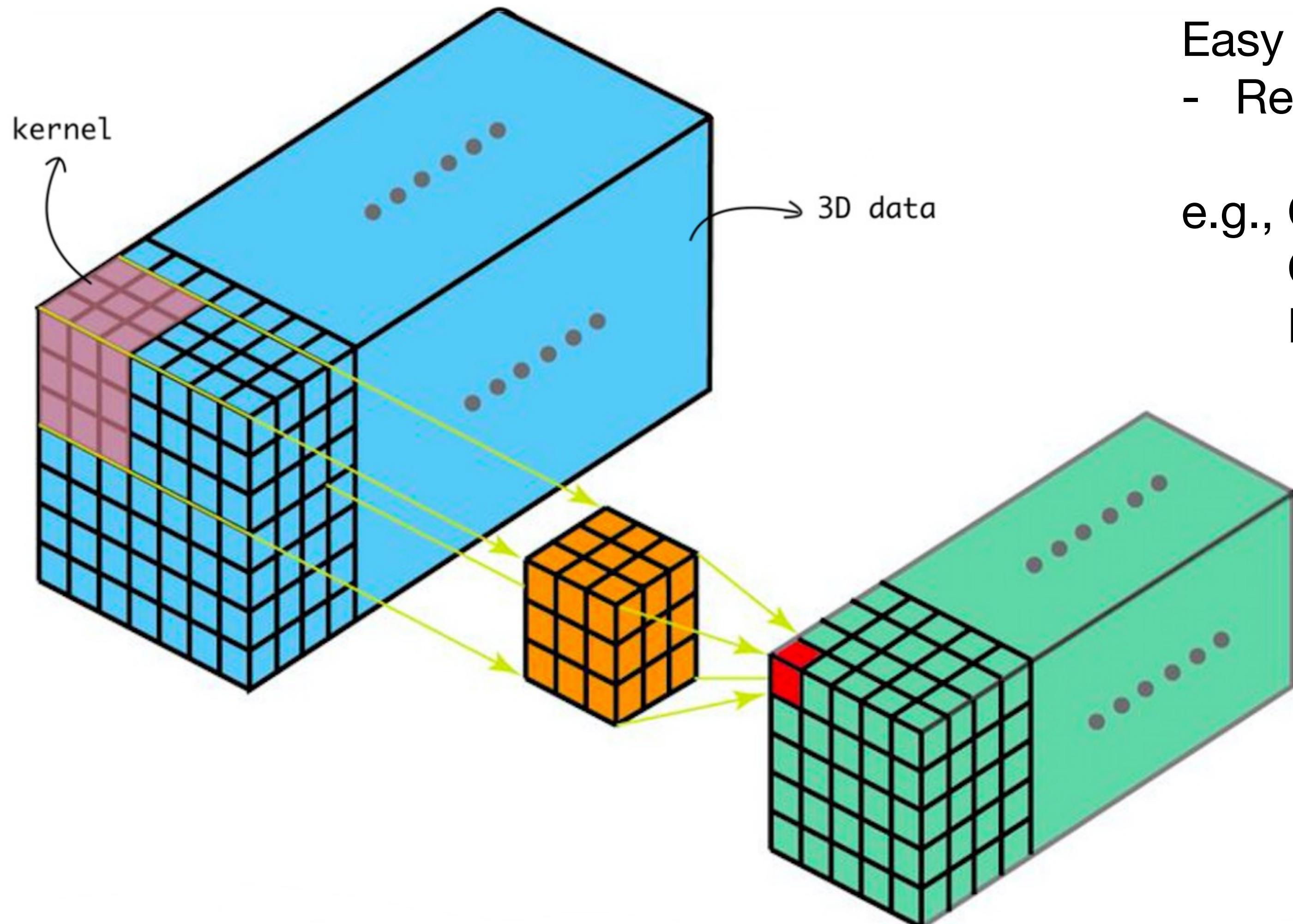


Toward 3D-aware Generative Models

3D Generative Adversarial Networks



3D Convolutional Layers



Easy to implement:

- Replace 2D by 3D in your code

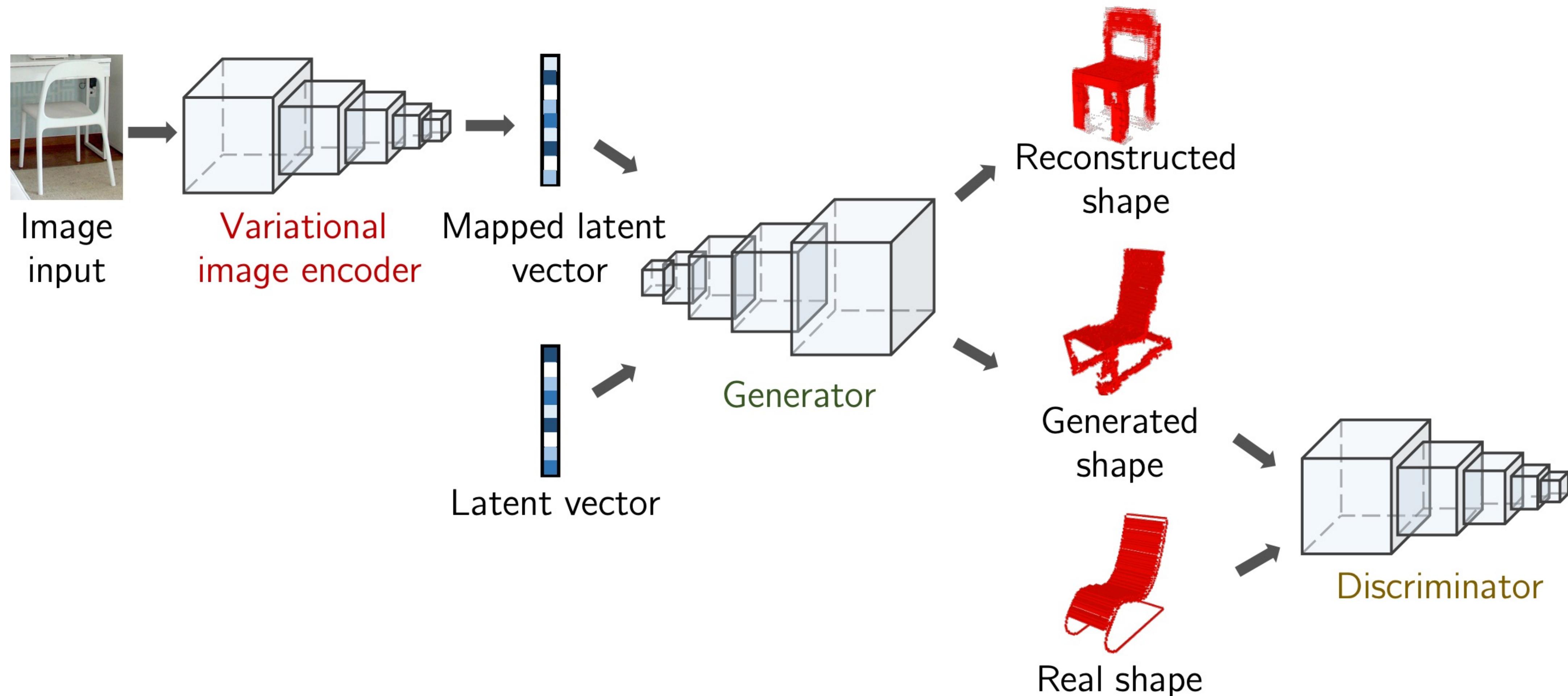
e.g., Conv2D -> Conv3D

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

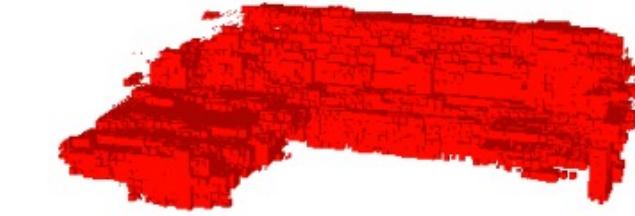


3D Generative Adversarial Networks



Input
image

Reconstructed
3D shape

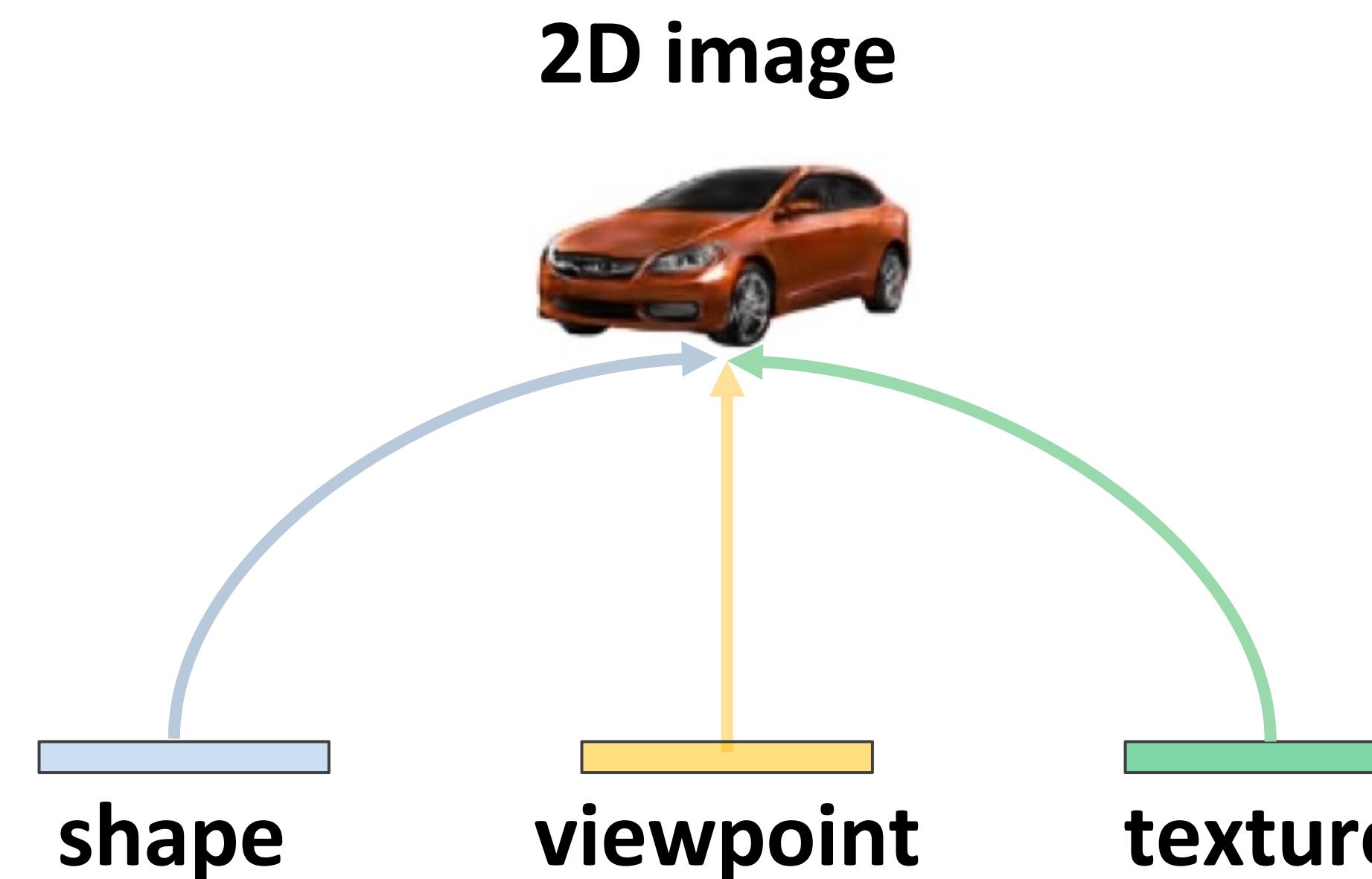


Input
image

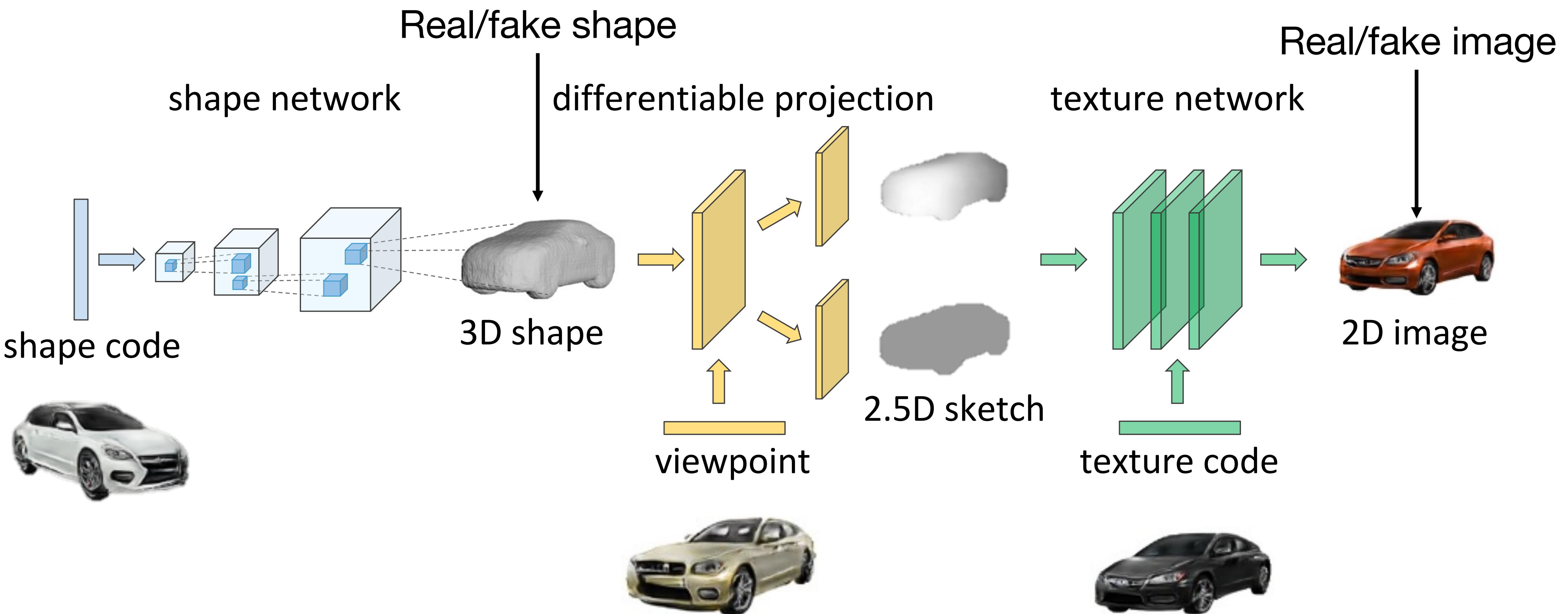
Reconstructed
3D shape

How to add Color and Texture?

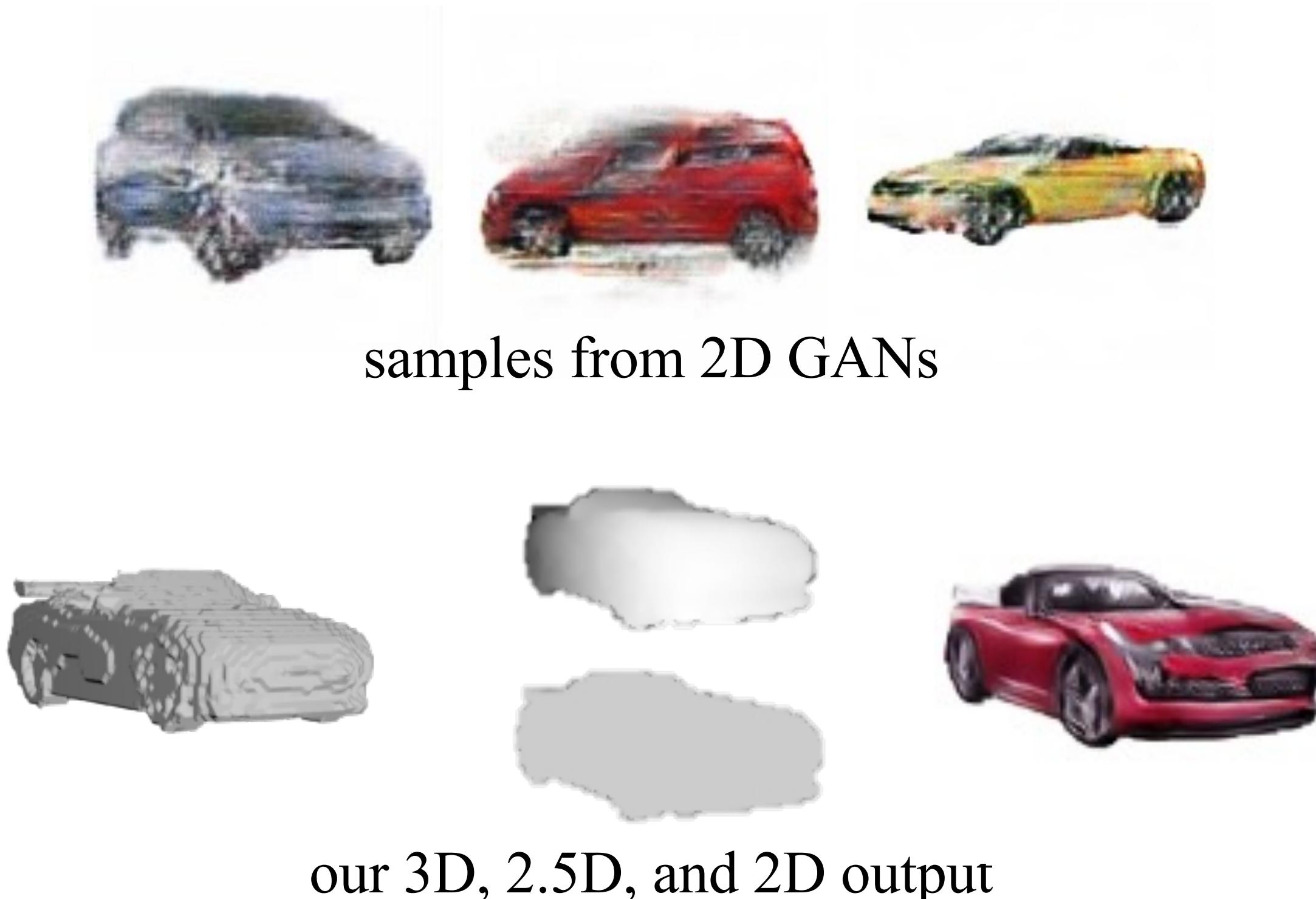
Learning 3D Disentanglement



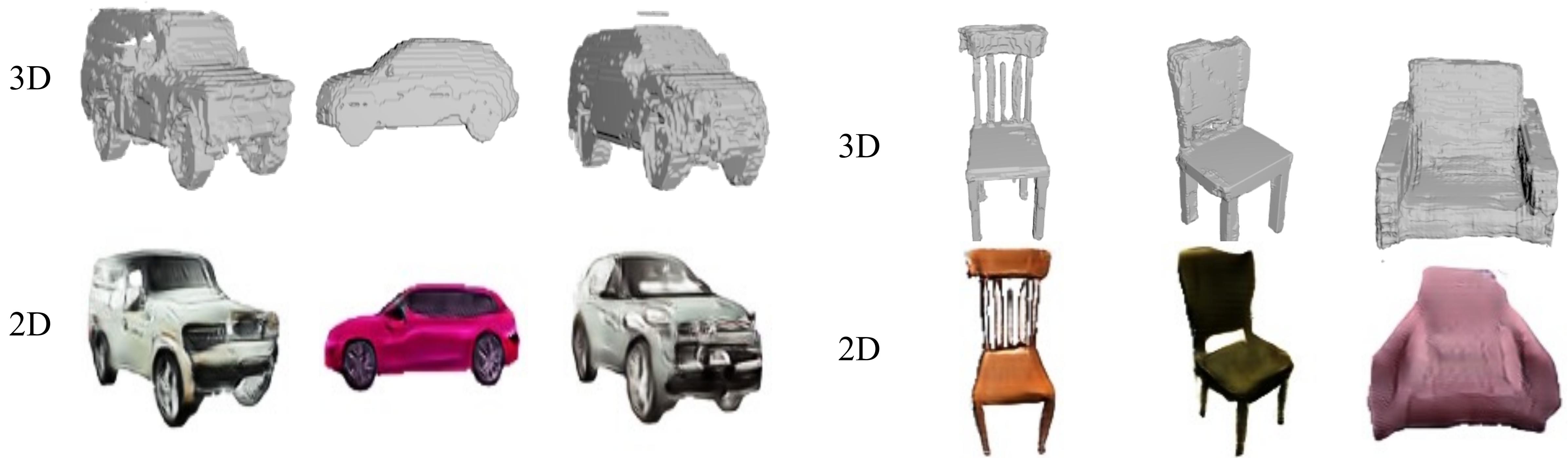
Learning 3D Disentanglement



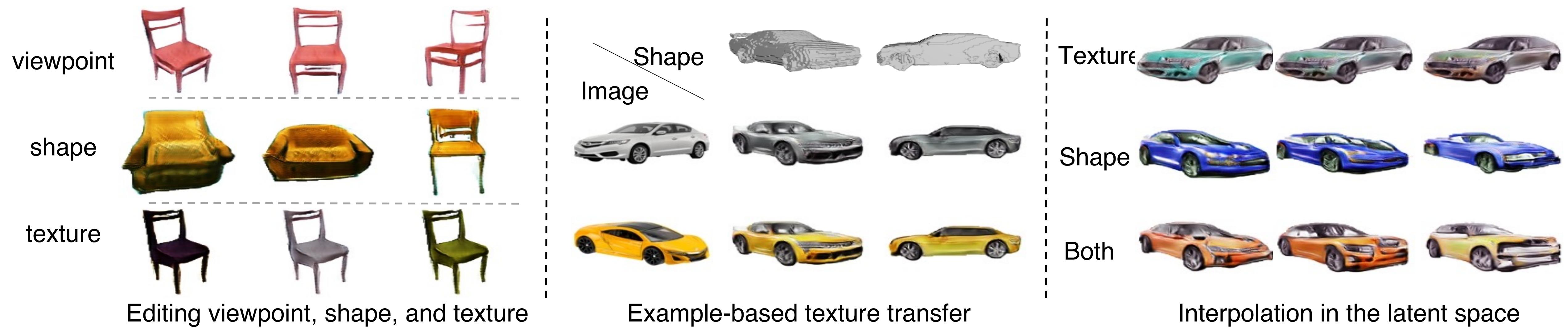
Learning 3D Disentanglement



Learning 3D Disentanglement



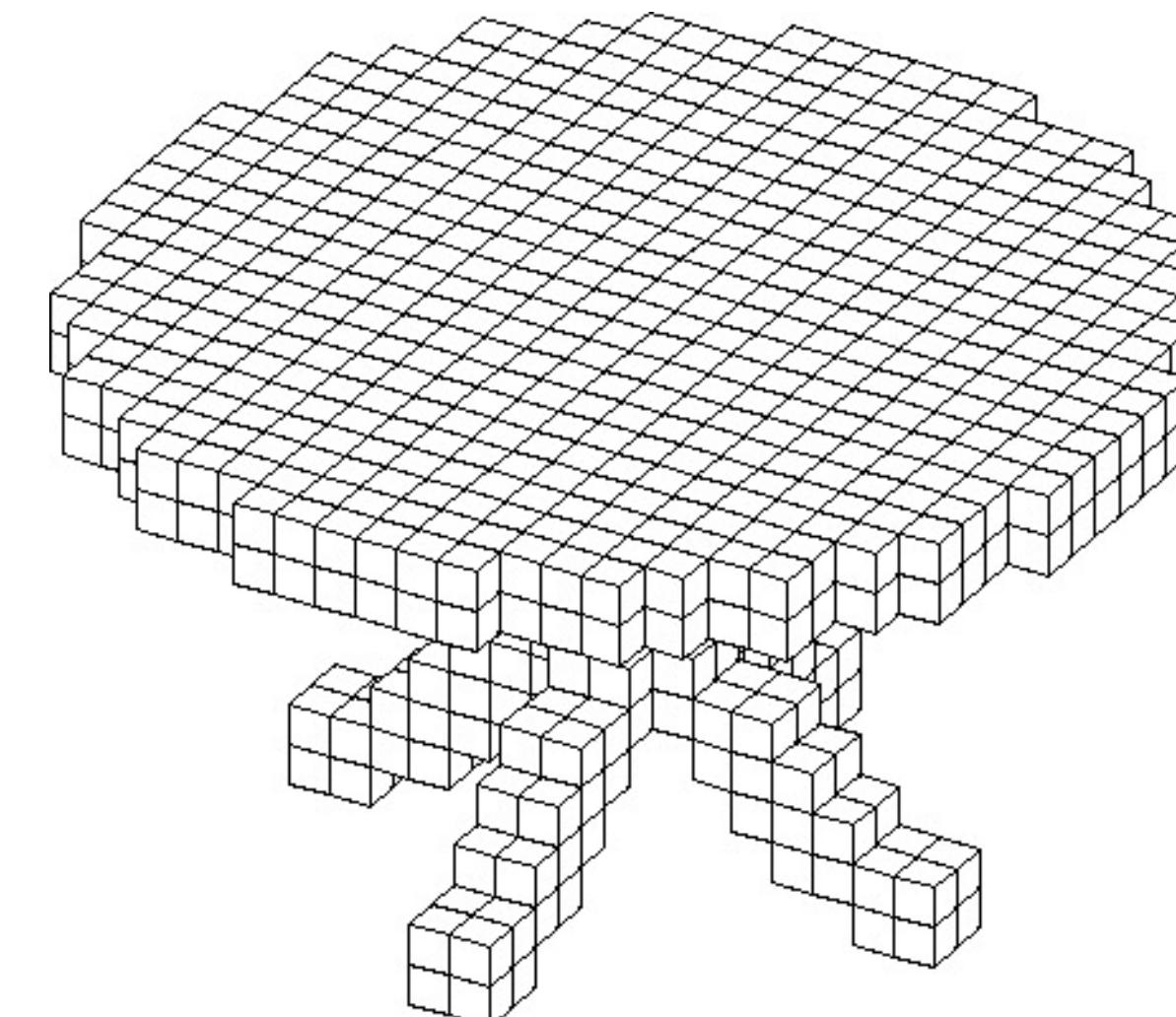
Learning 3D Disentanglement



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

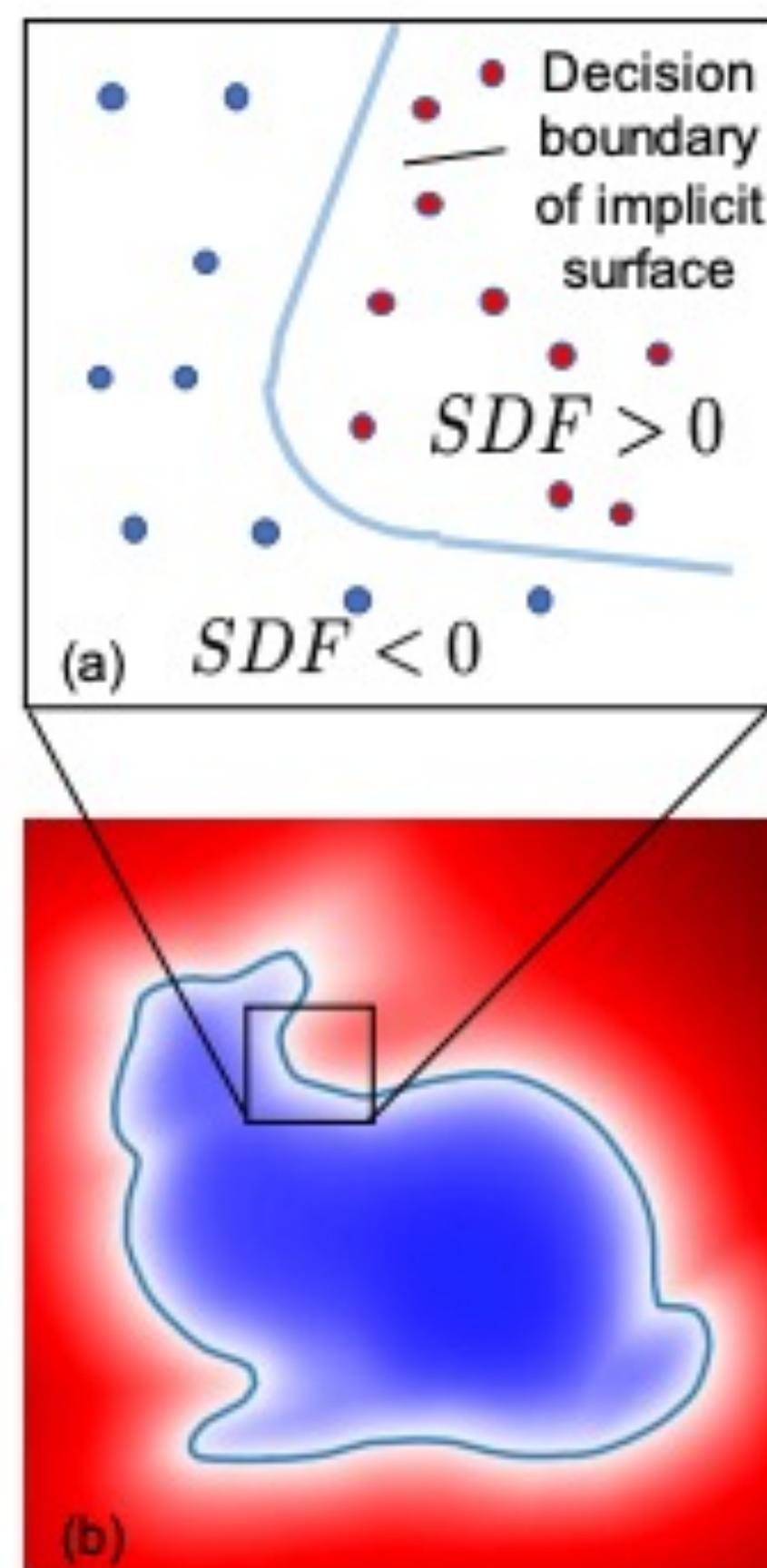
256x256x256 -> 1024x1024x1024

Cannot even fit a single training data to GPU

Improvements:

1. Using implicit representation (network-based)

Signed Distance Function (SDF)



Explicit function:

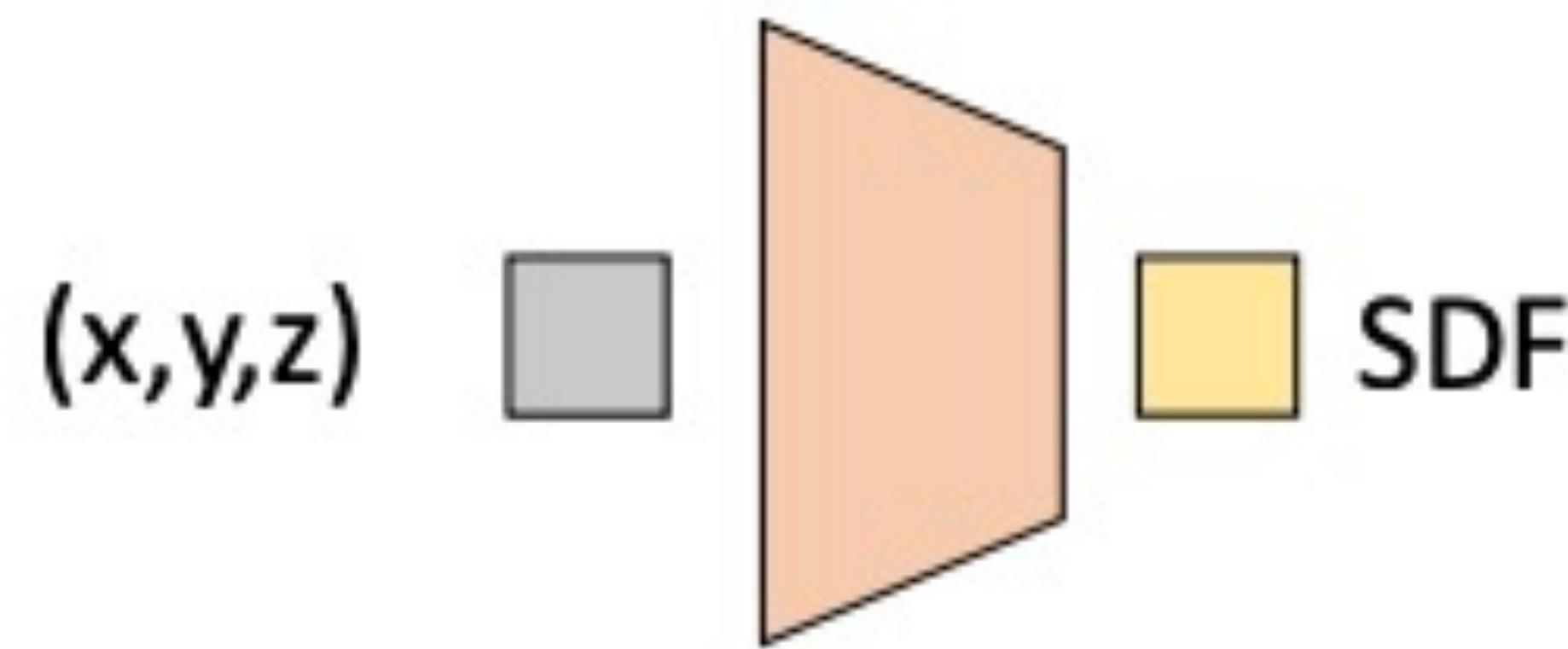
$$y = 2x. \quad (y = f(x))$$

Implicit function:

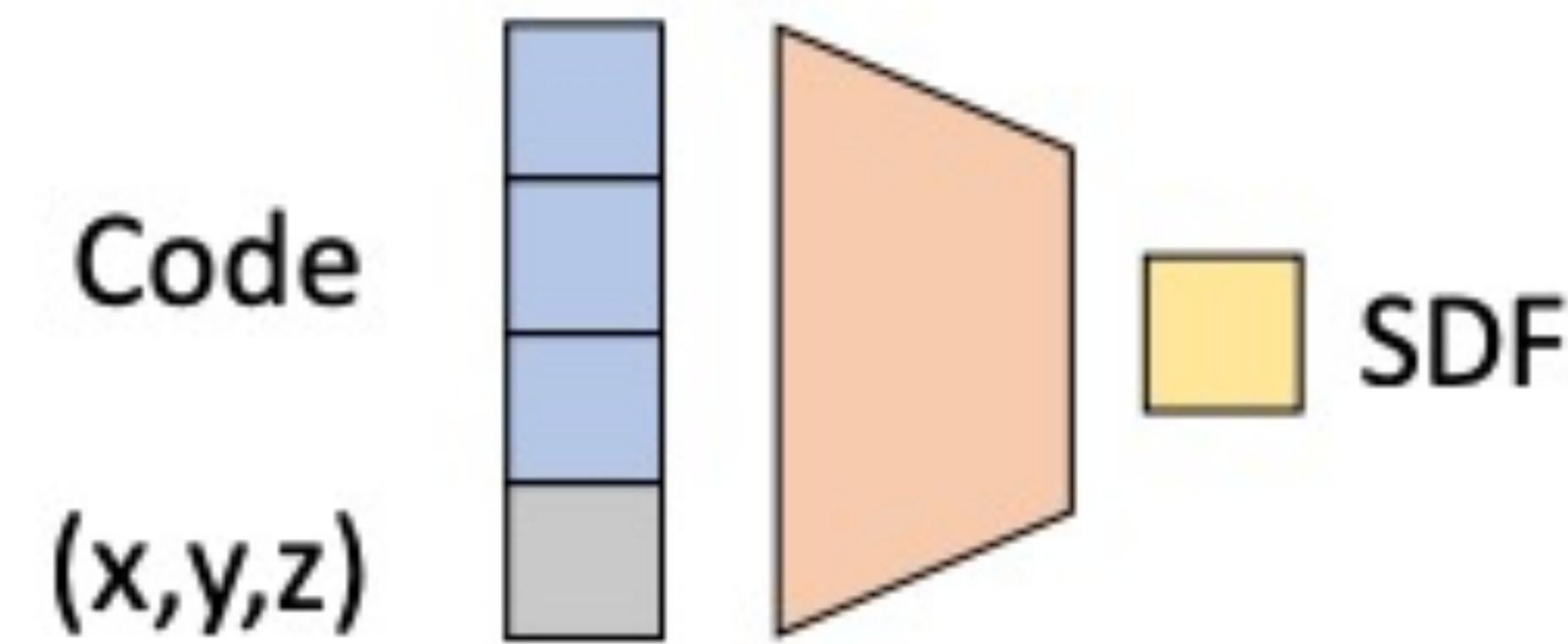
$$2y - 4x = 0, F(x, y) = 0$$

A set of zeros of a function of two variables.

Deep SDF

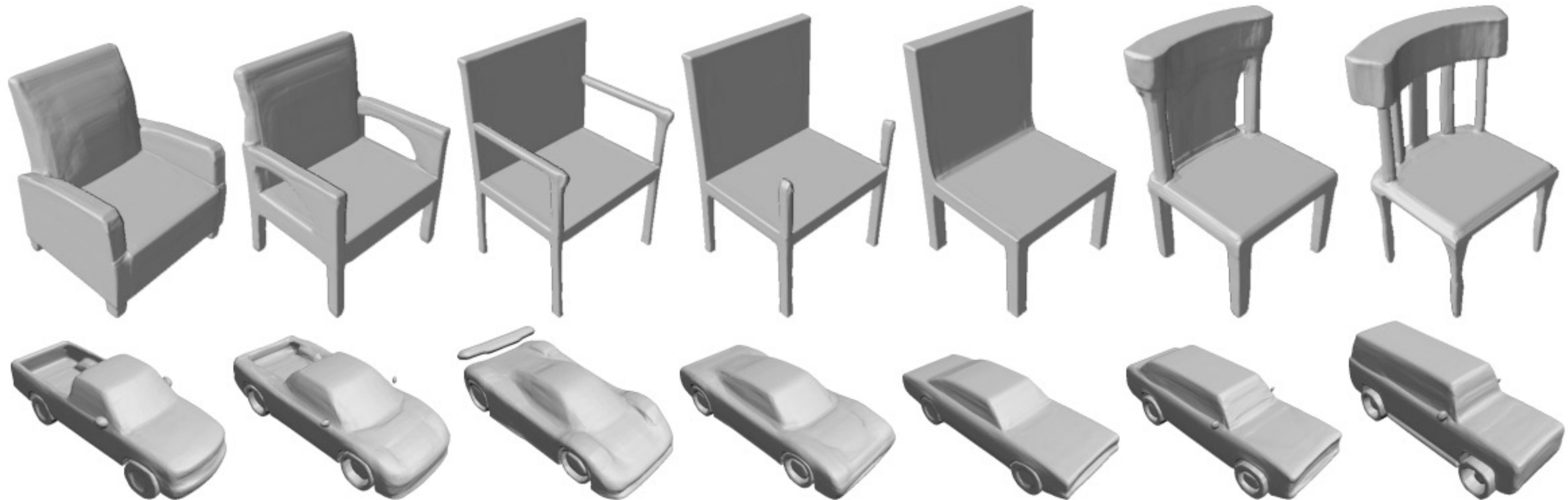


(a) Single Shape DeepSDF

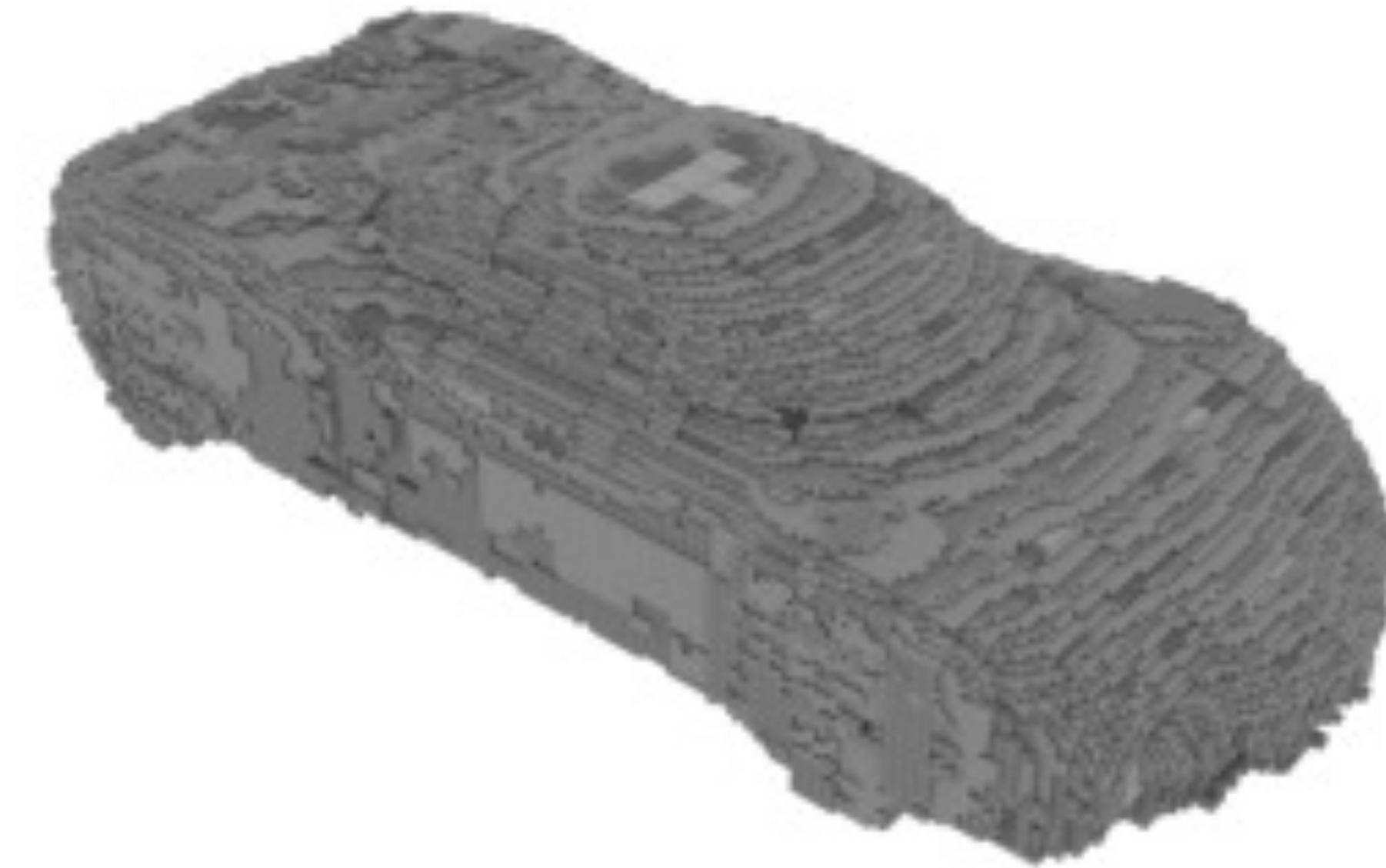


(b) Coded Shape DeepSDF

Deep SDF



Deep SDF

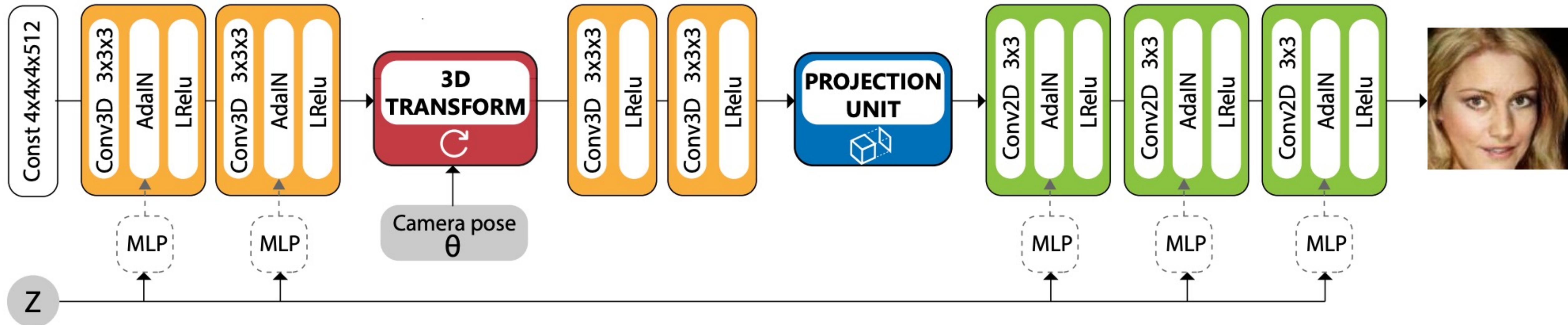


DeepSDF preserve details and render visually pleasing results compared to voxel-based methods.

Improvements:

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

HoloGAN

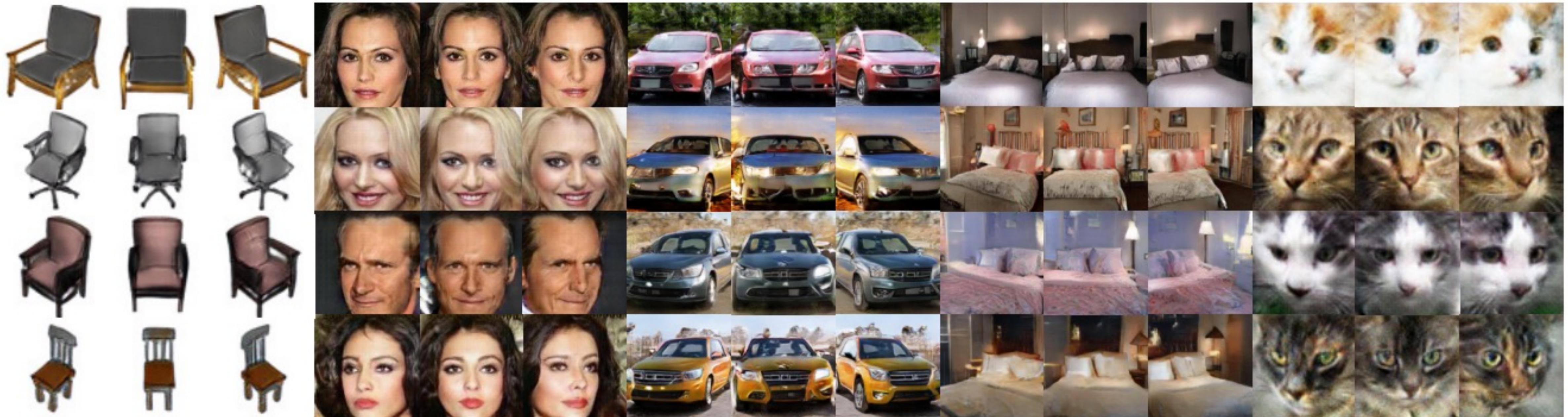


Representation: 3D feature representation

Training: Adversarial loss + latent code reconstruction

Modulation: AdaIN

HoloGAN



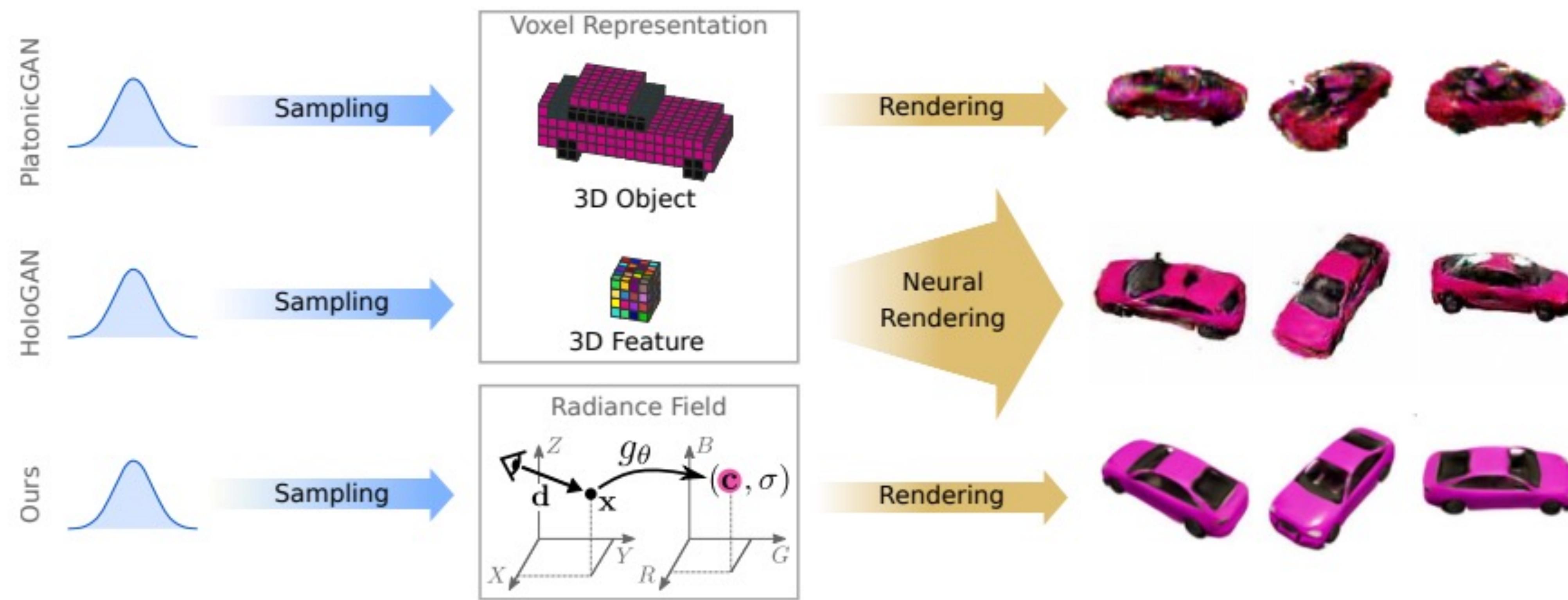
Limitations:

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

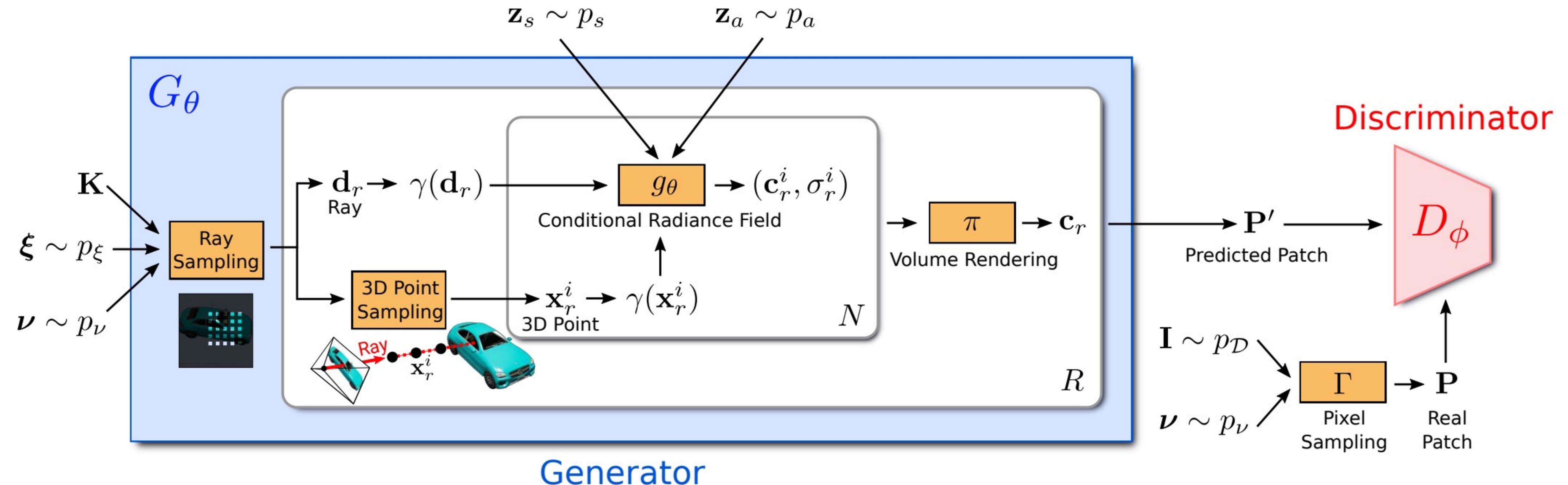
NeRF + GANs

(Neural rendering + Generative Models)

GRAF: Generative Radiance Fields

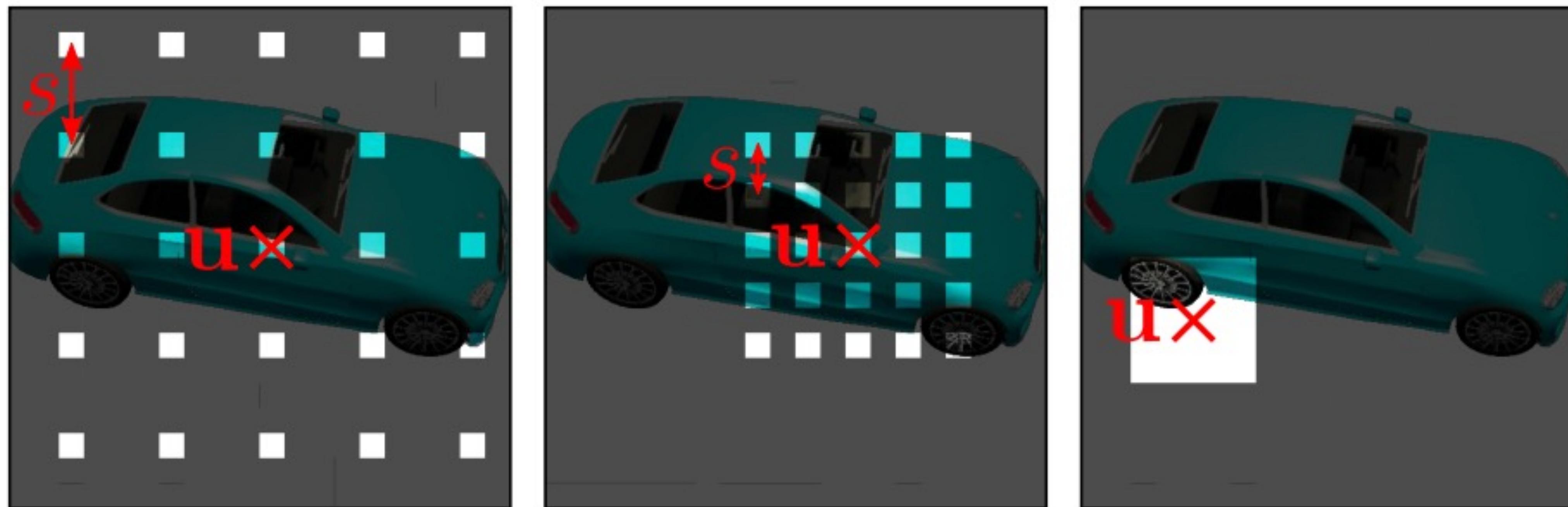


GRAF: Generative Radiance Fields



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

GRAF: Generative Radiance Fields

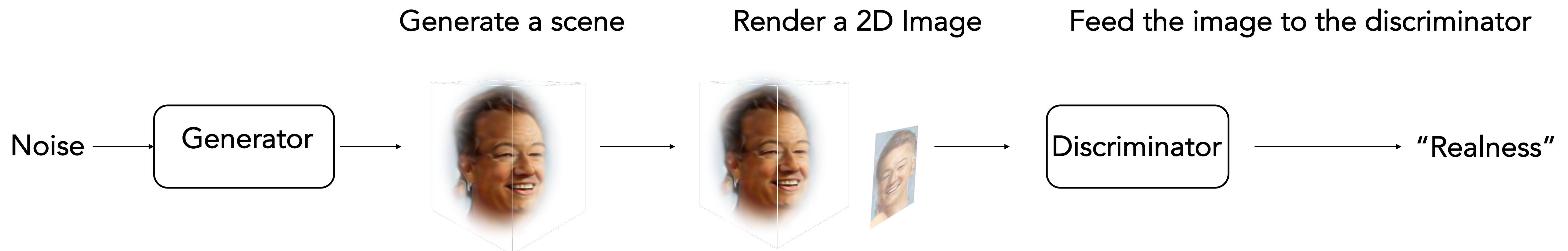


Multi-scale ray sampling

Training a 3D-Aware GAN

3D-Aware GAN Training Steps

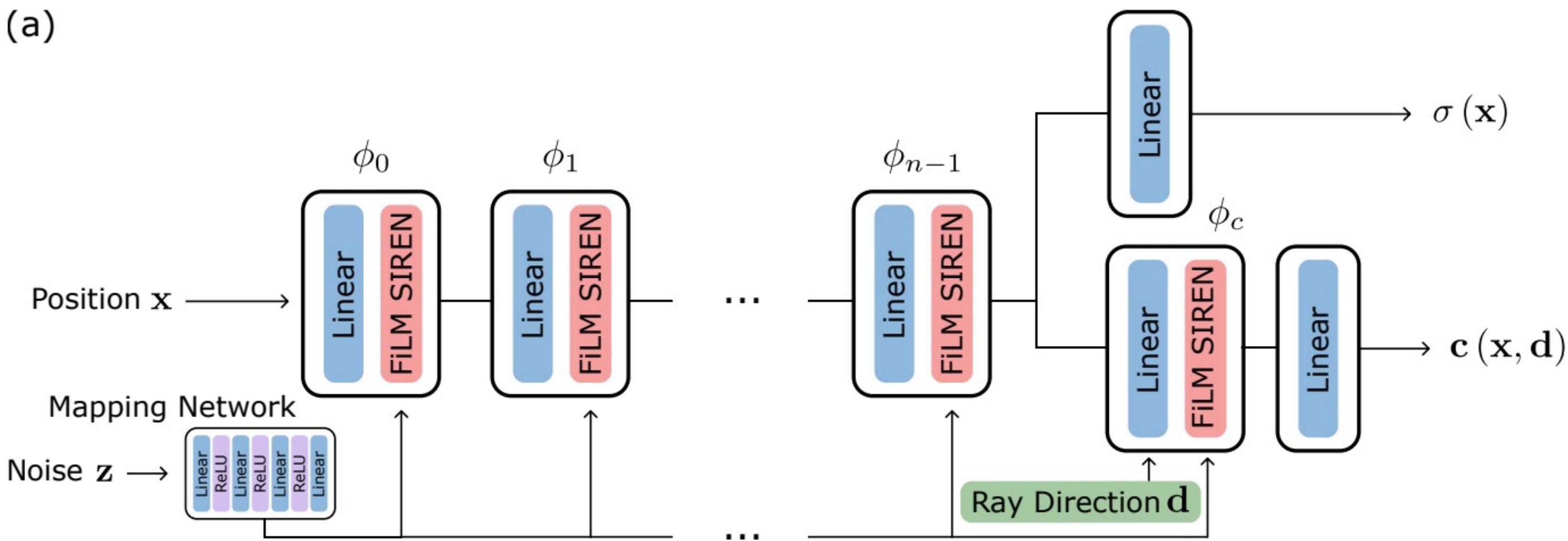
1. Generate a representation of a scene
2. Render the scene from a random camera pose
3. Feed the image to a 2D discriminator
4. Backpropagate through the discriminator and differentiable rendering



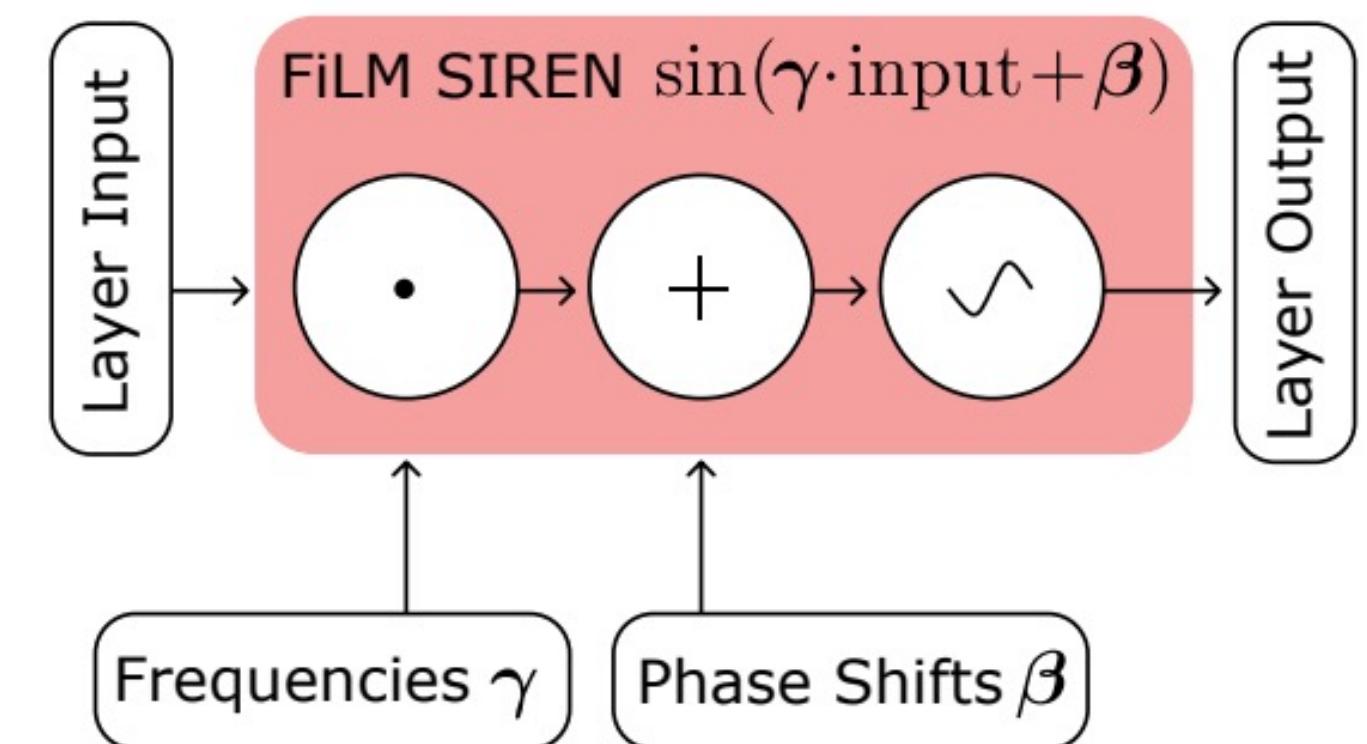
Slide credit: Eric Chan

π -GAN

(a)



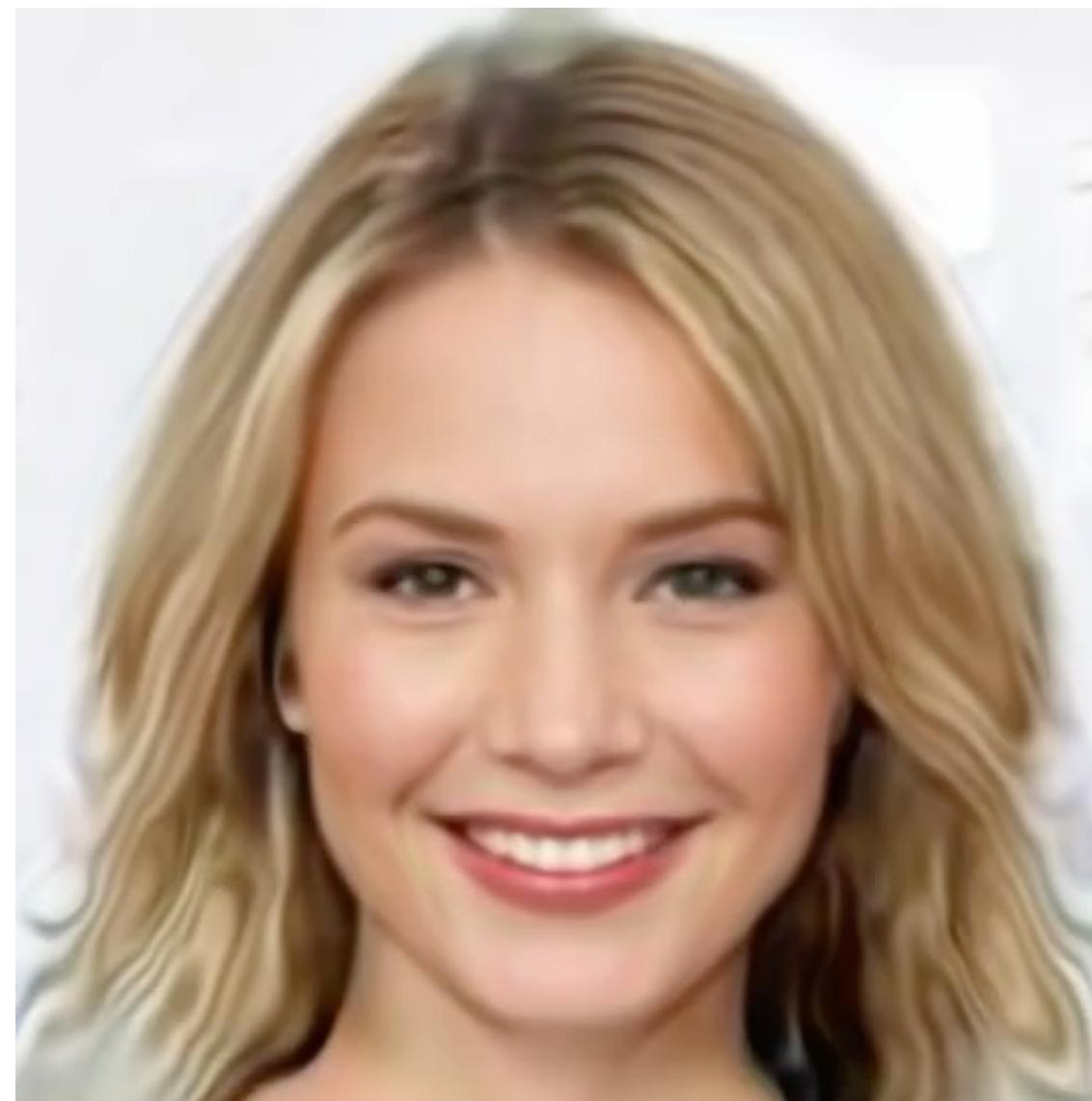
(b)



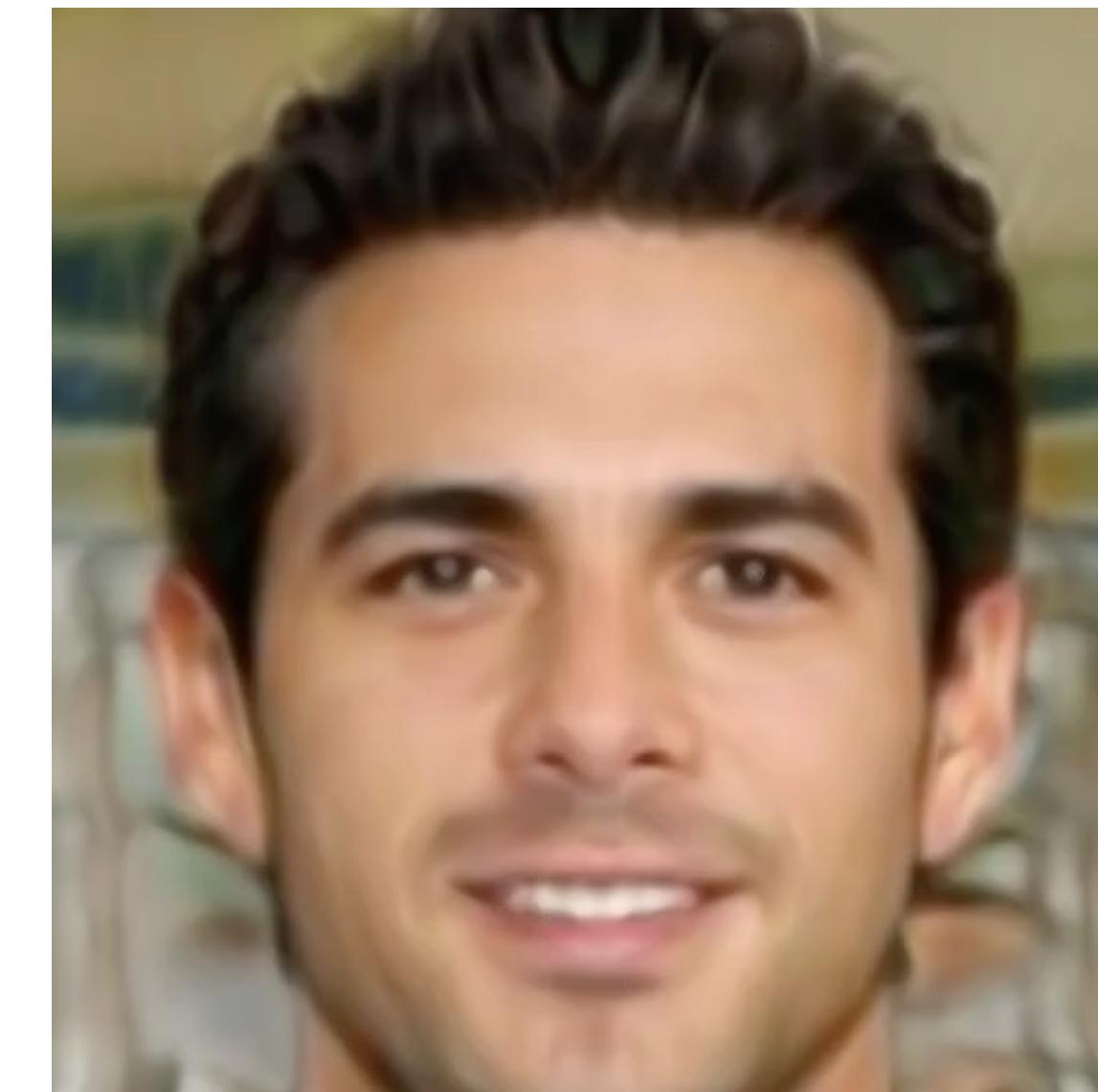
Mapping network + AdalN (FiLM) + learnable positional encoding

$$\phi_i(\mathbf{x}_i) = \sin(\boldsymbol{\gamma}_i \cdot (\mathbf{W}_i \mathbf{x}_i + \mathbf{b}_i) + \boldsymbol{\beta}_i)$$

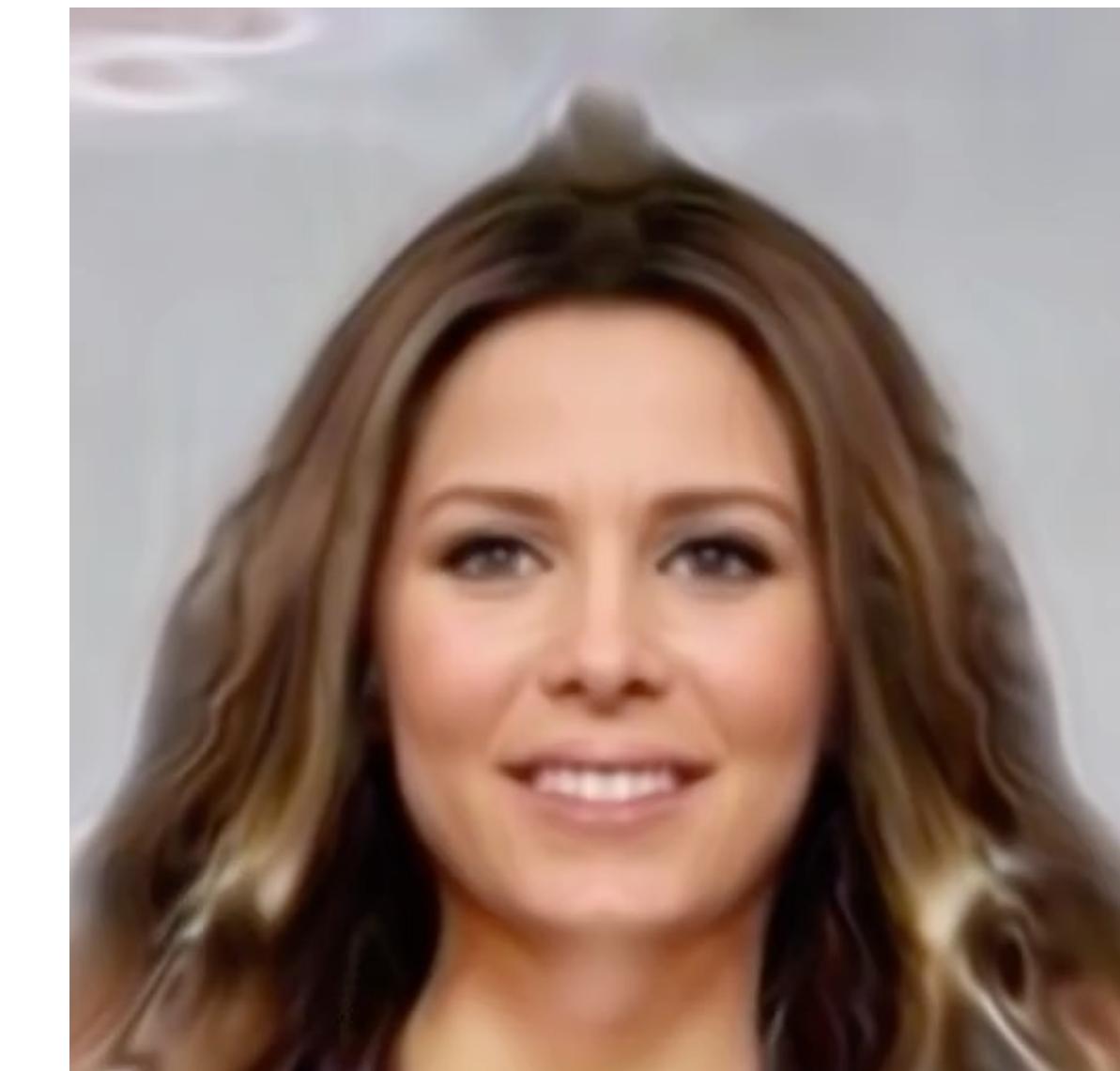
π -GAN



Focal Length



Camera Position



Latent Interpolation

Slide credit: Eric Chan

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

π -GAN

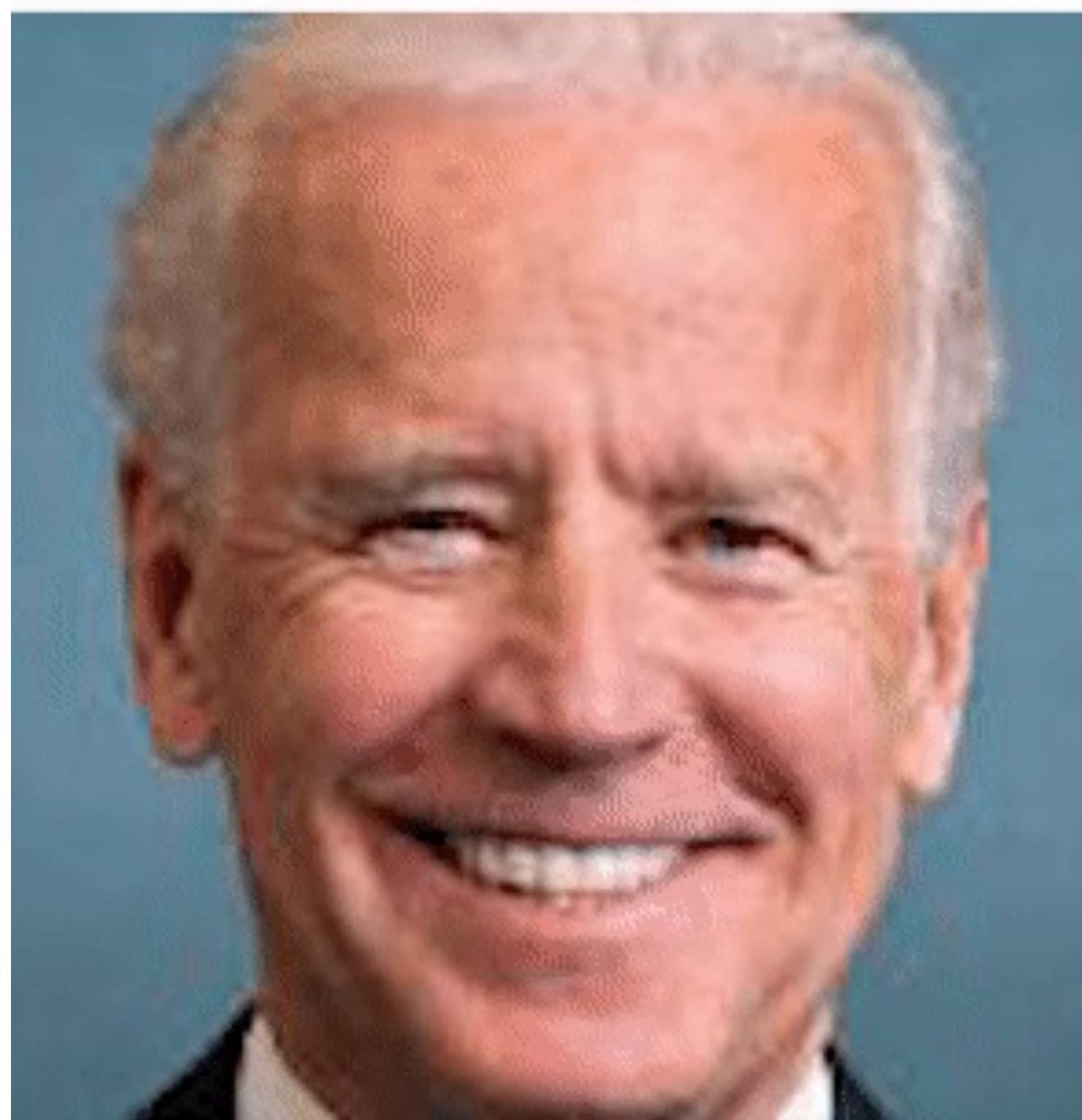


Slide credit: Eric Chan

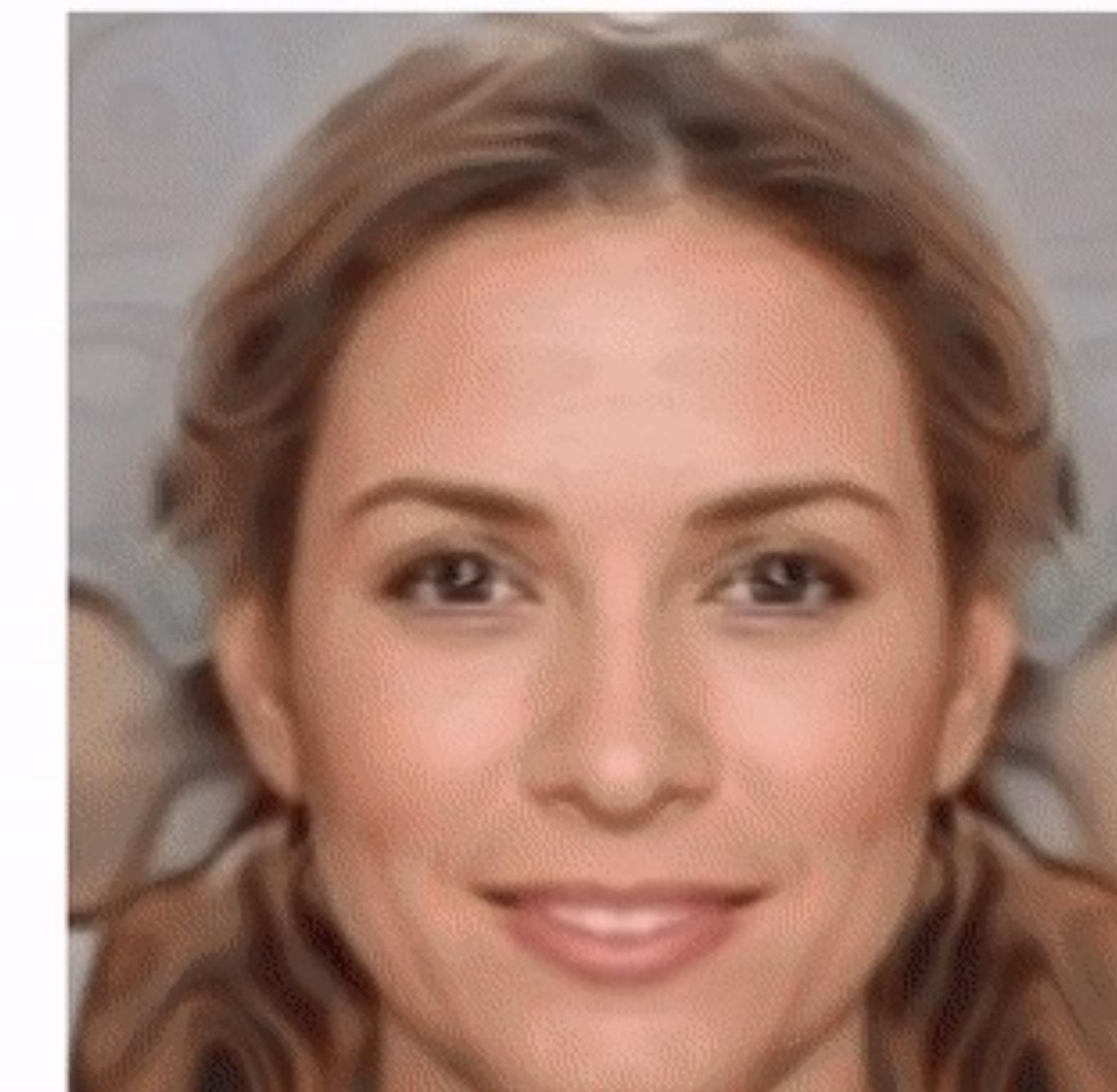
pi-GAN: Periodic Implicit Generative Adversarial Networks for 3D-Aware Image Synthesis [Chan et al., 2021]

π -GAN

Target



Reconstruction



Slide credit: Eric Chan

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

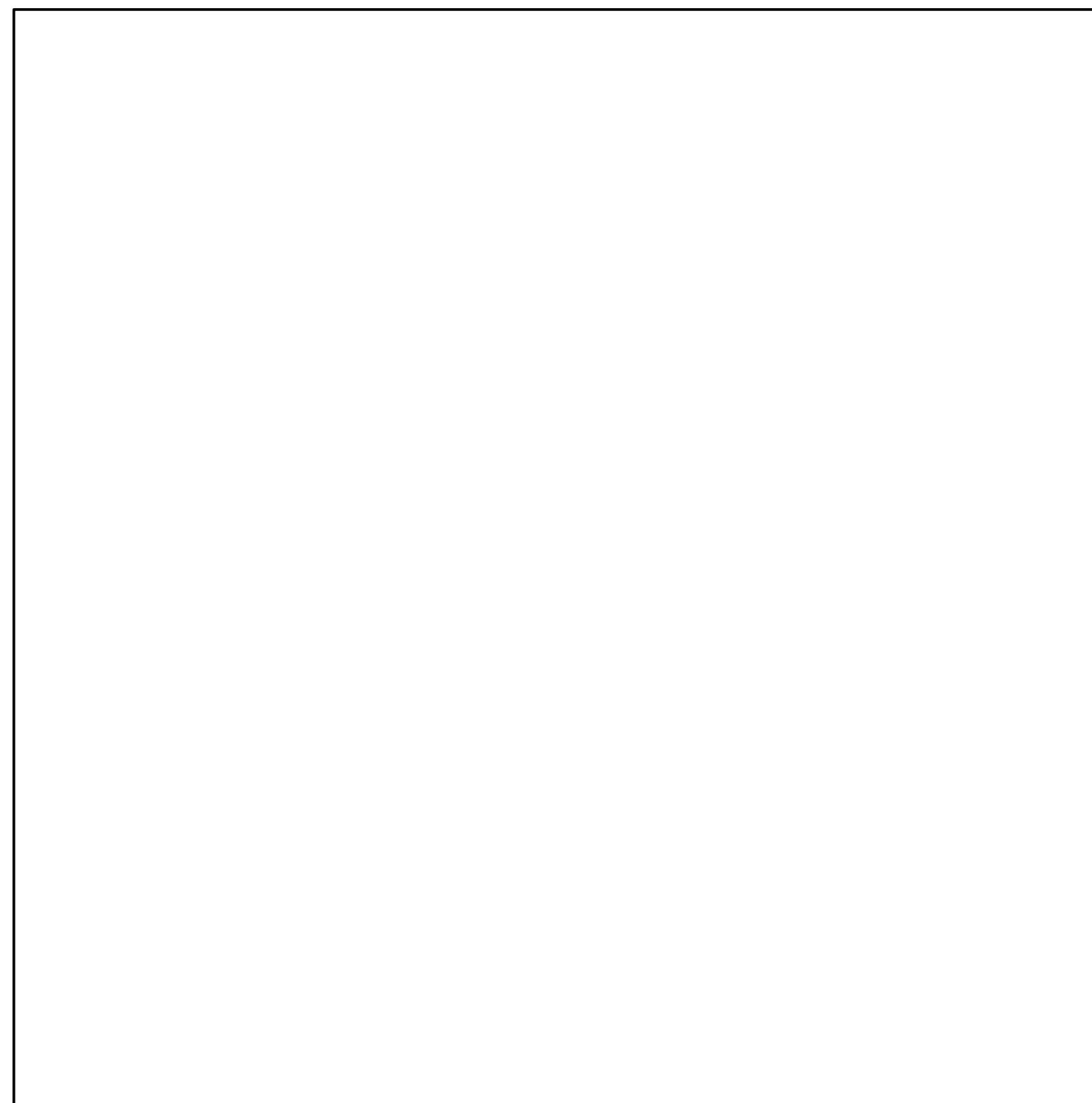
3D Object Editing



Input Views

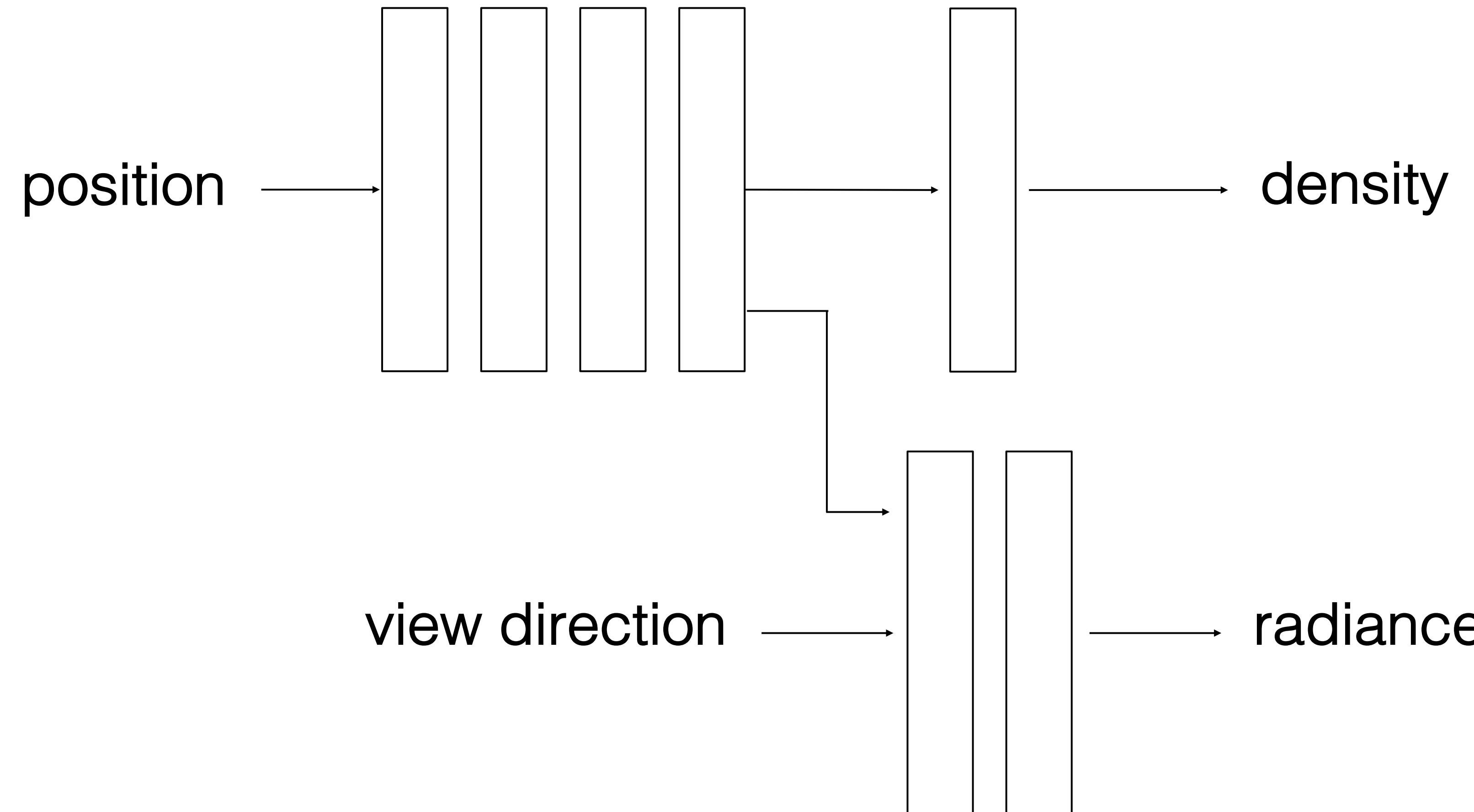


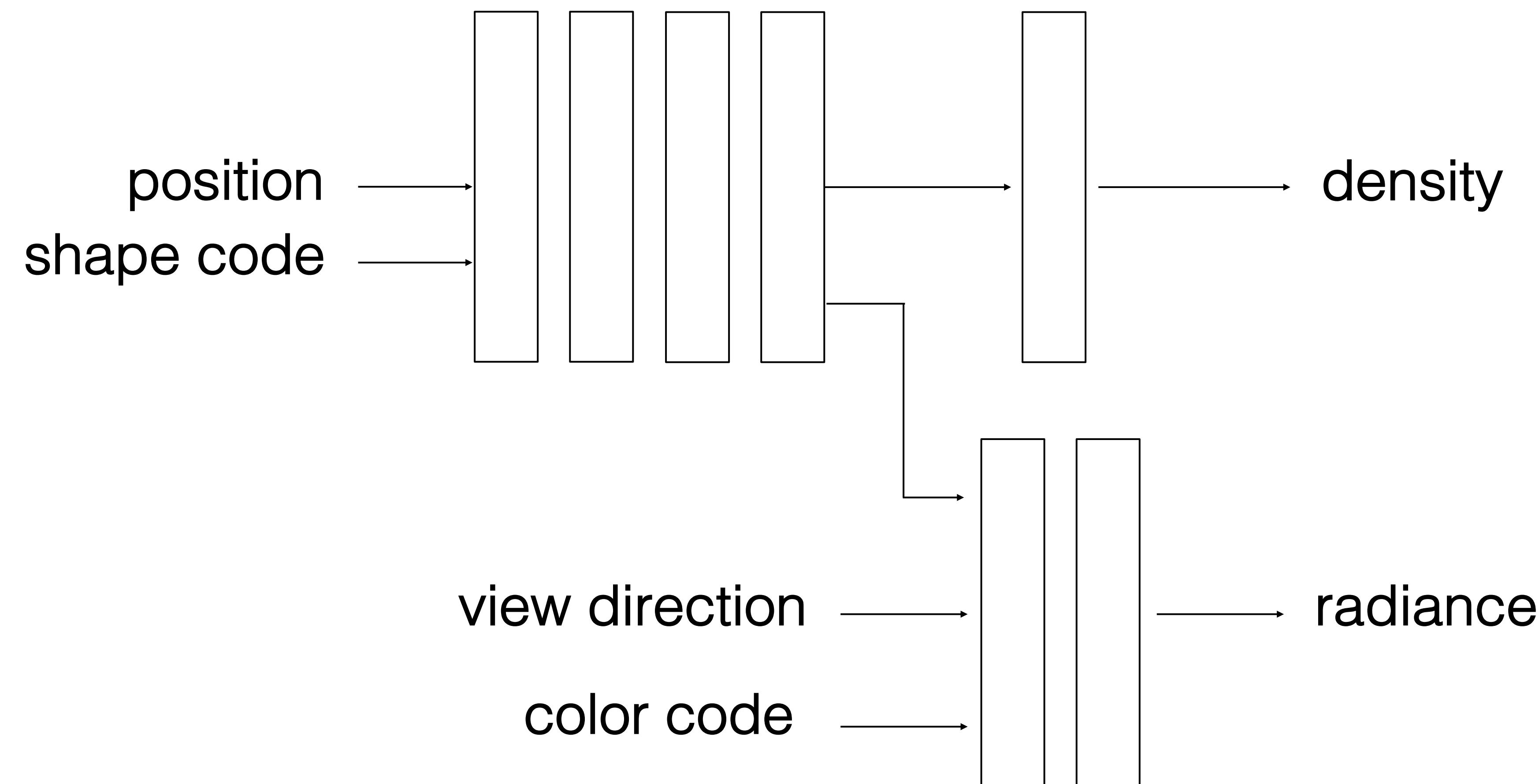
Output Novel Rendered Views

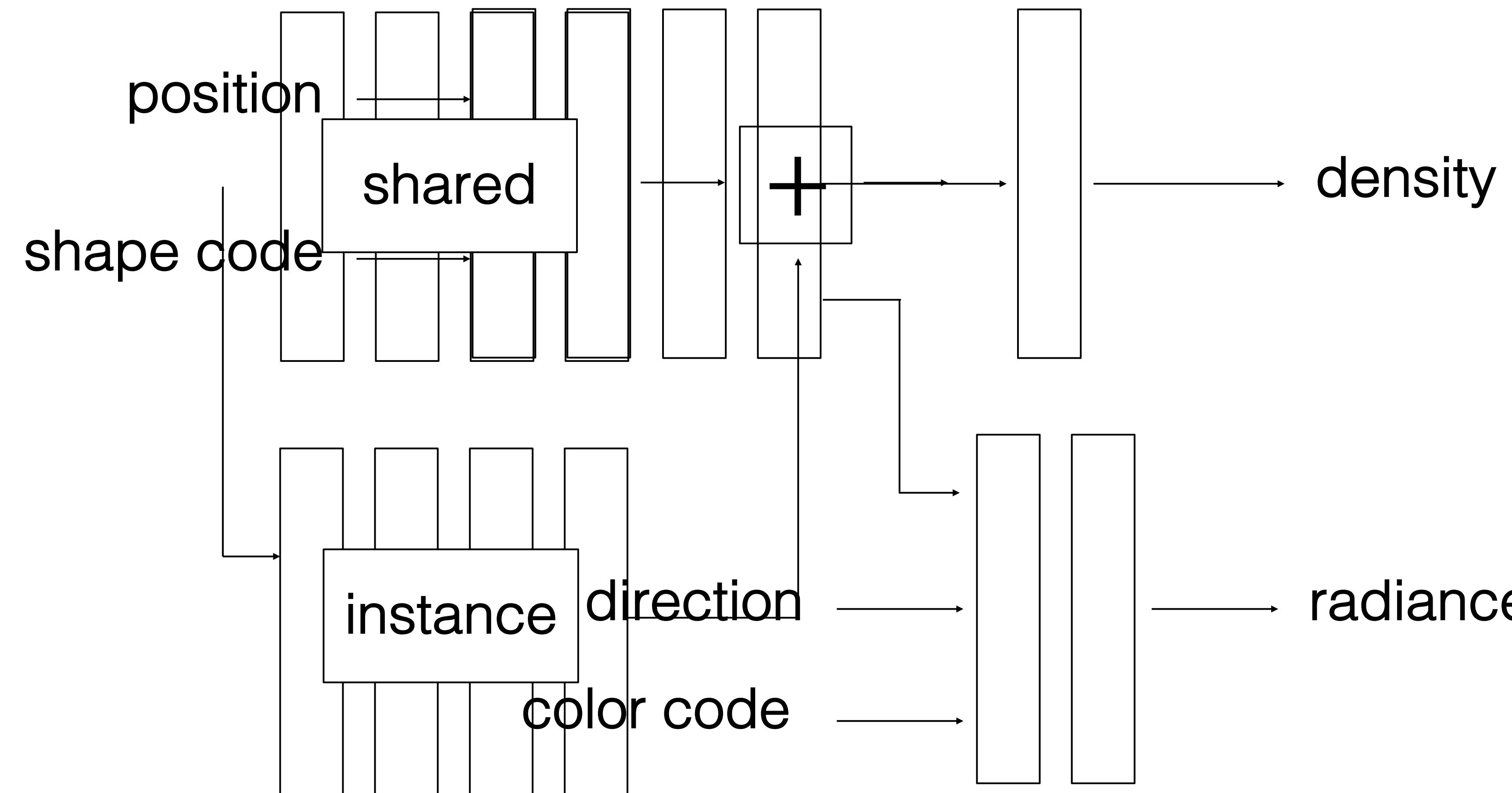


Editing Conditional Radiance Fields [Liu et al., 2021]

Neural Radiance Fields Base Architecture



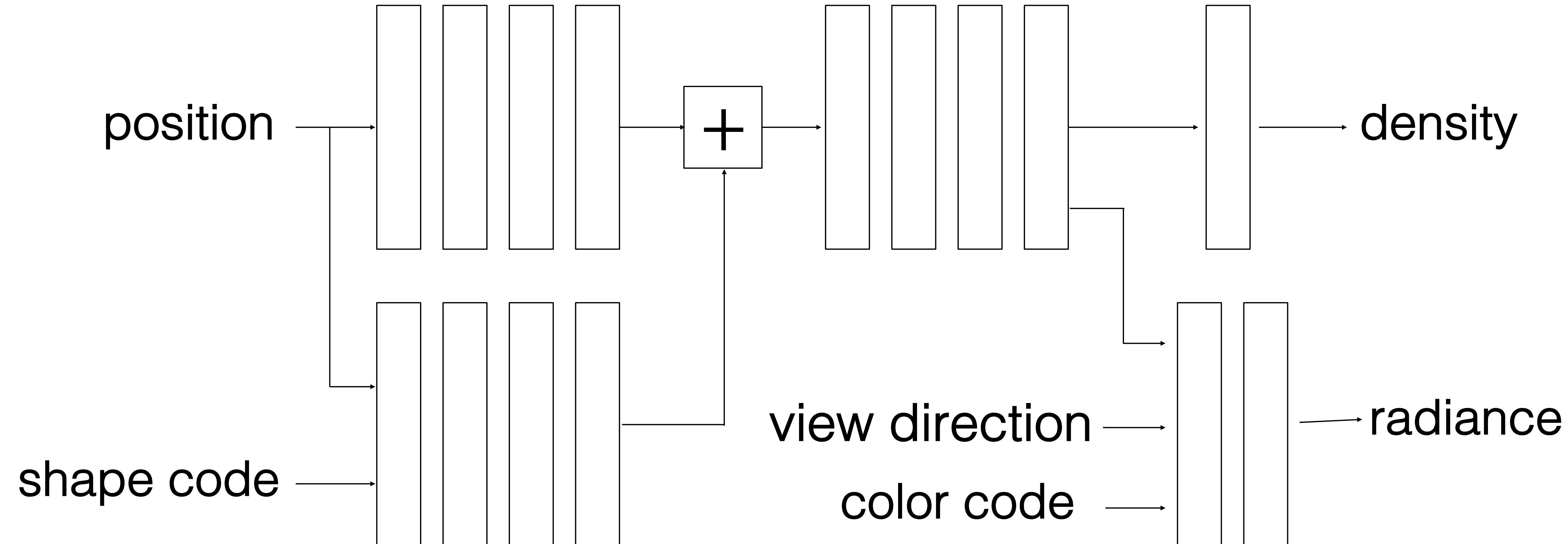




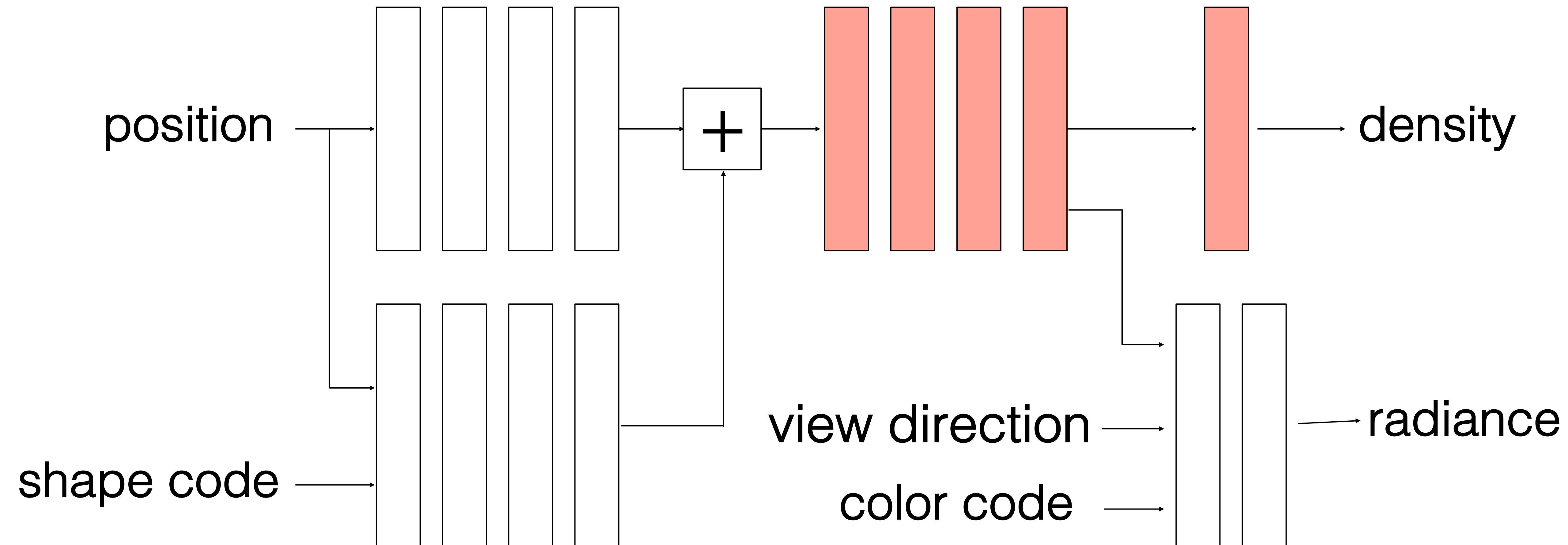


Editing Conditional Radiance Fields [Liu et al., 2021]

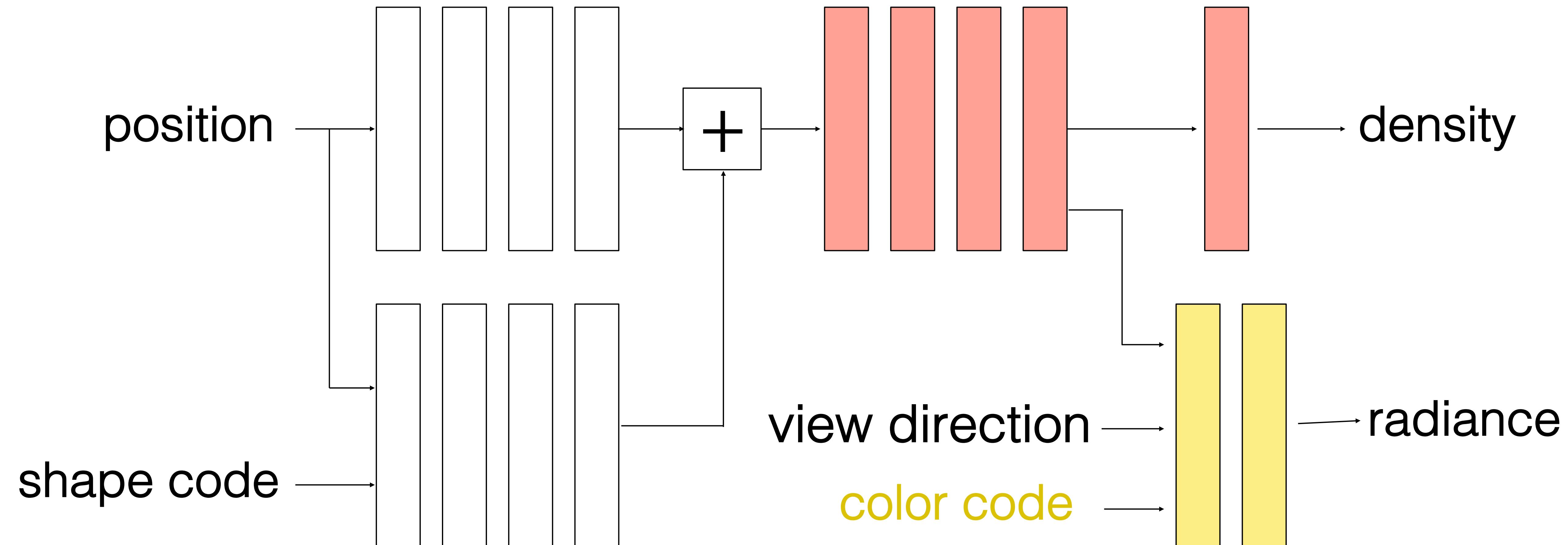
Which parameters
do we change?



Updated for Shape Editing



Updated for Shape Editing



Updated for Color Editing

Color Editing



Input User Scribble

Output Edited Views

Shape Editing



Input User Scribble



Output Edited Views

Color Editing



Input User Scribble



Output Edited Views

Shape Editing

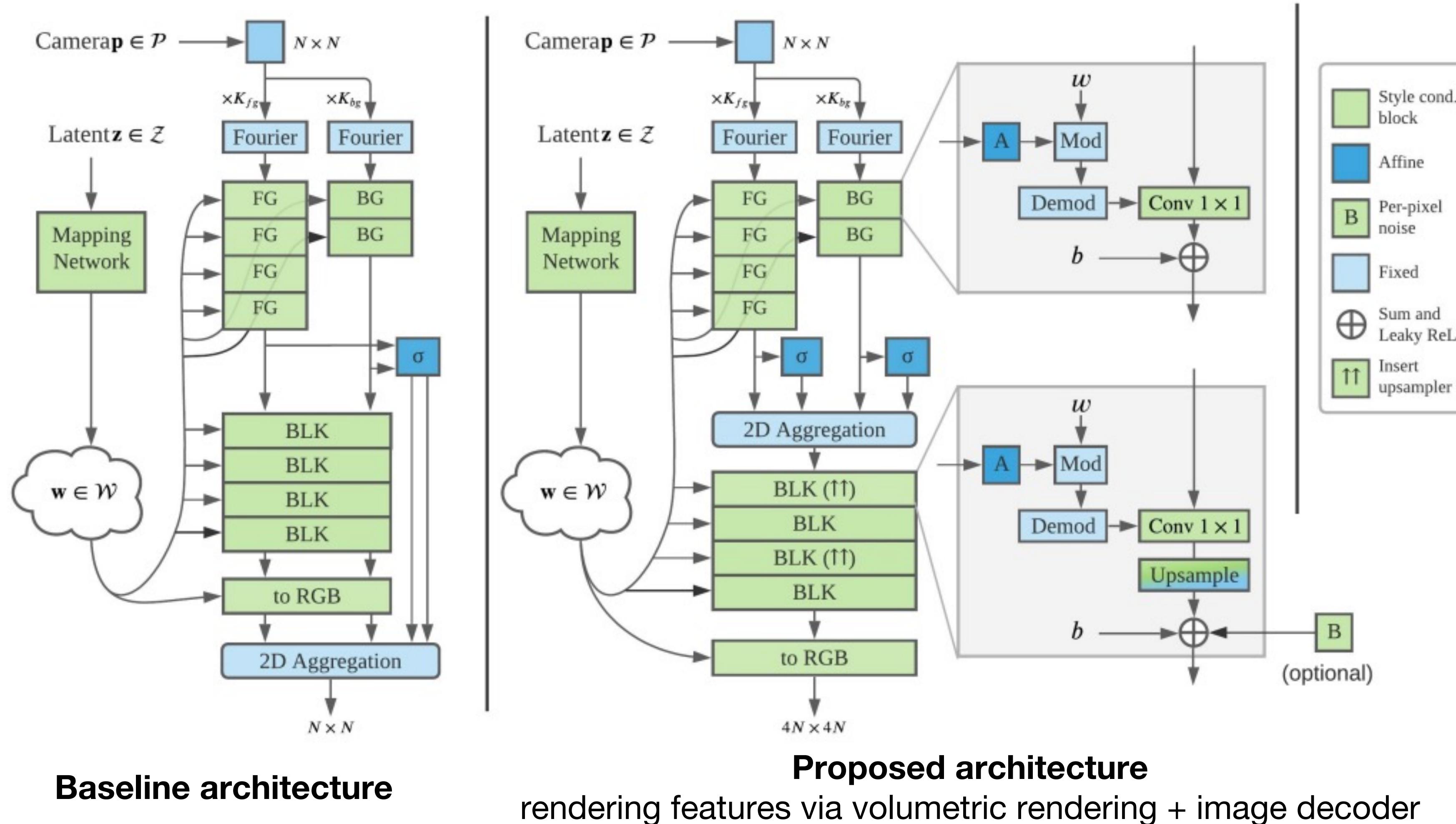


Input User Scribble



Output Edited Views

Advanced Architectures: StyleNeRF



Baseline architecture

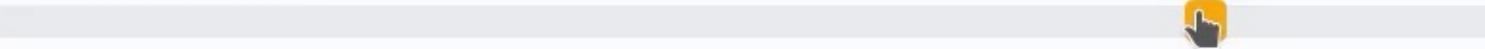
Proposed architecture
rendering features via volumetric rendering + image decoder

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

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

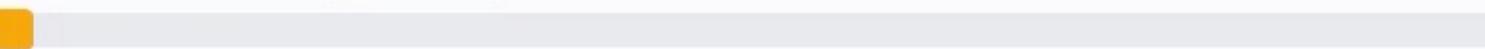
MODEL NAME
FFHQ512 ▾

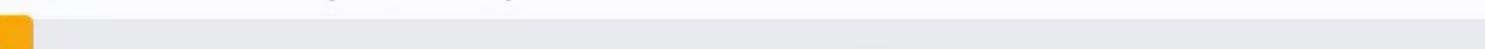
CHECKPOINT PATH

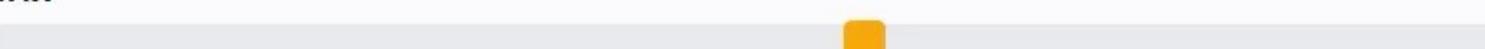
TRUNCATION TRICK  0.7

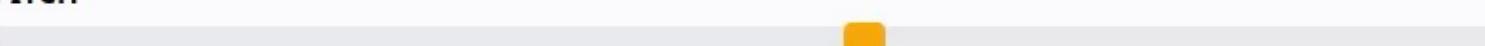
SEED1
4

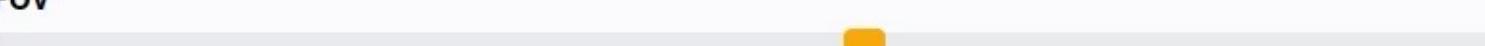
SEED2
9

LINEAR MIXING RATIO (GEOMETRY)  0

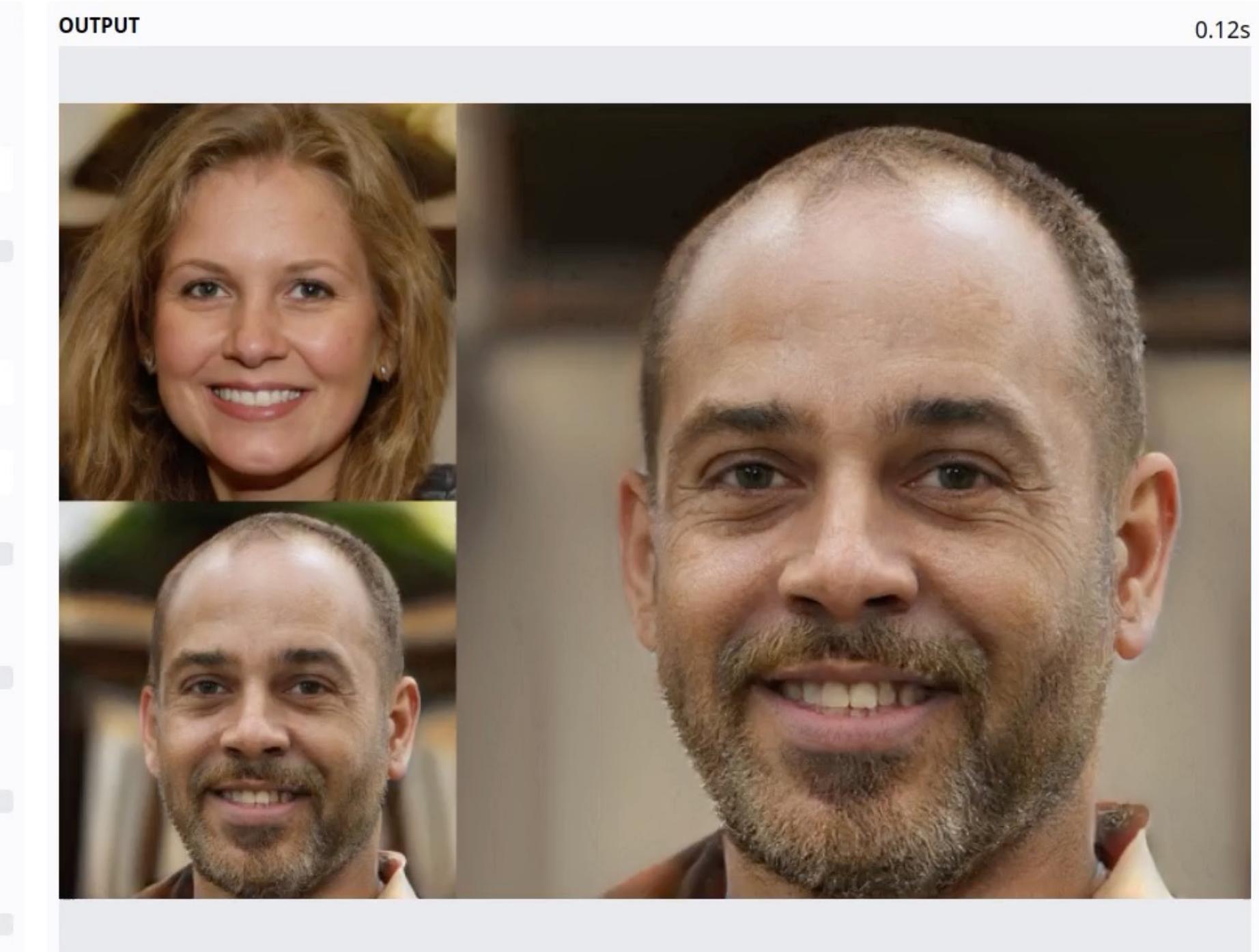
LINEAR MIXING RATIO (APPARENCE)  0

YAW  0

PITCH  0

FOV  12

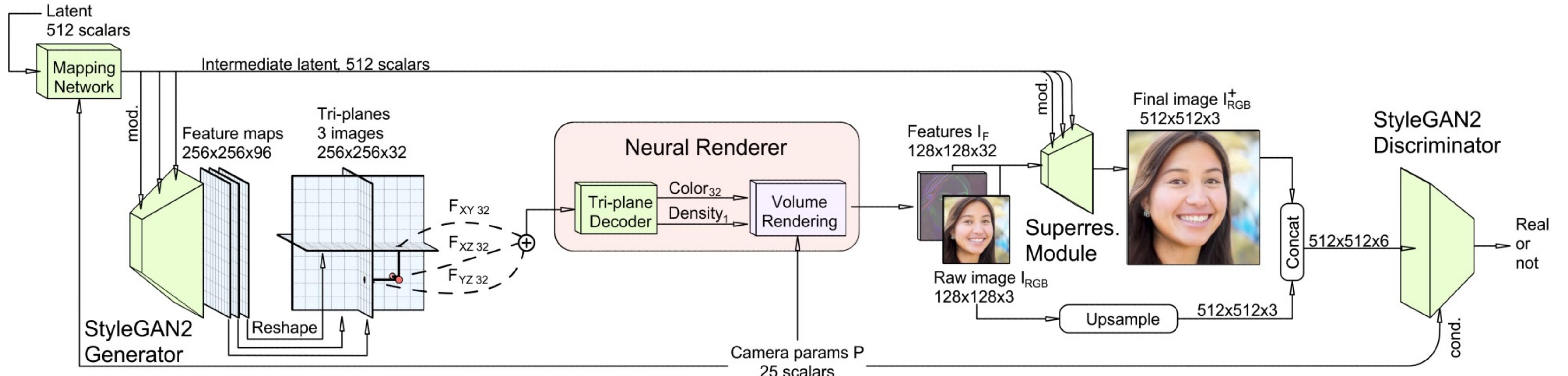
Clear



Screenshot

Flag

Advanced Architectures: StyleNeRF



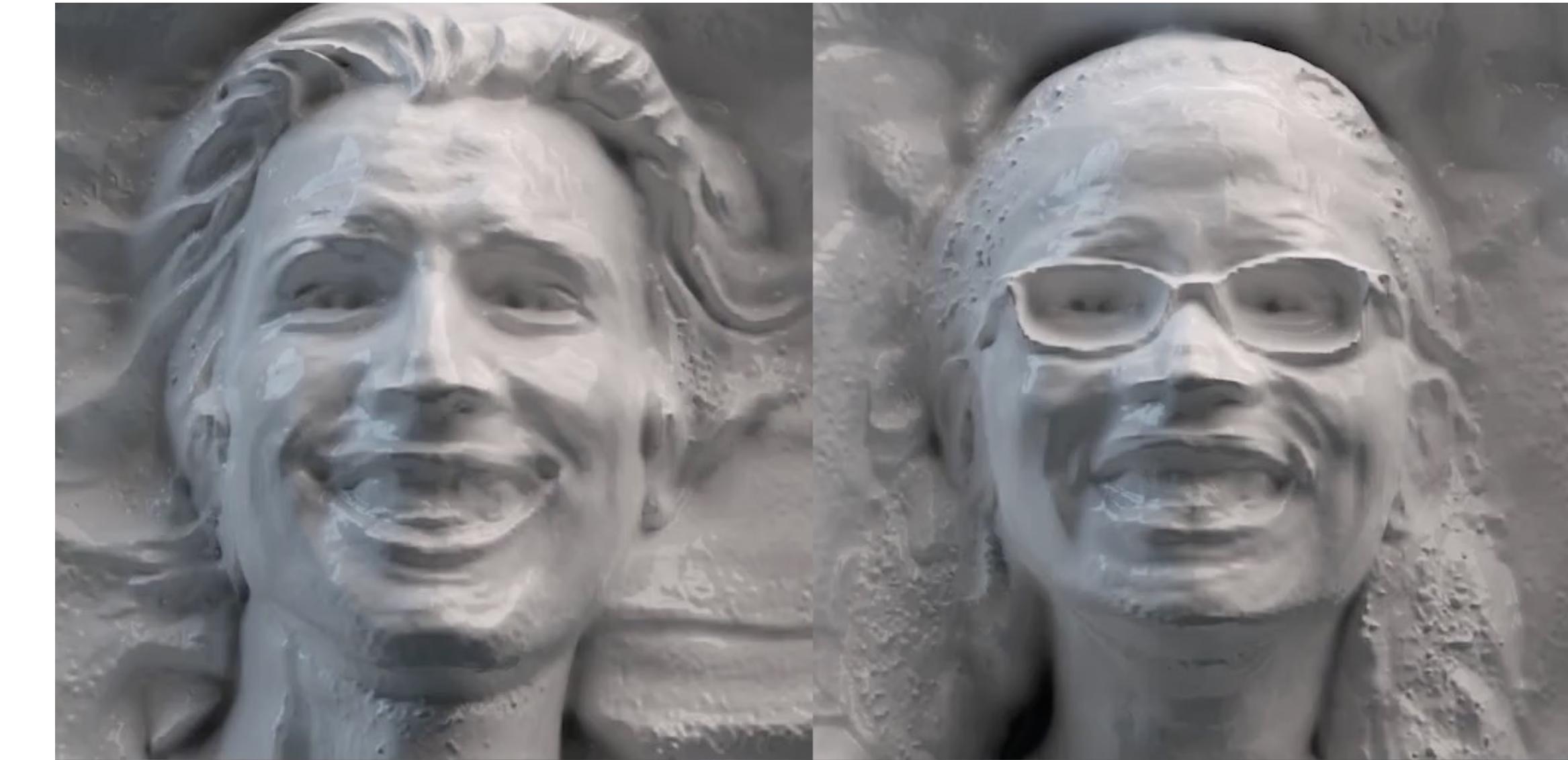
**Tri-plane representation
for speed-up**

**Rendering features
via volumetric rendering**

**Generate final output
via image encoder**

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

Thank You!

<https://learning-image-synthesis.github.io/sp22/>