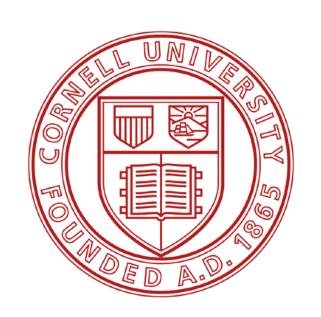
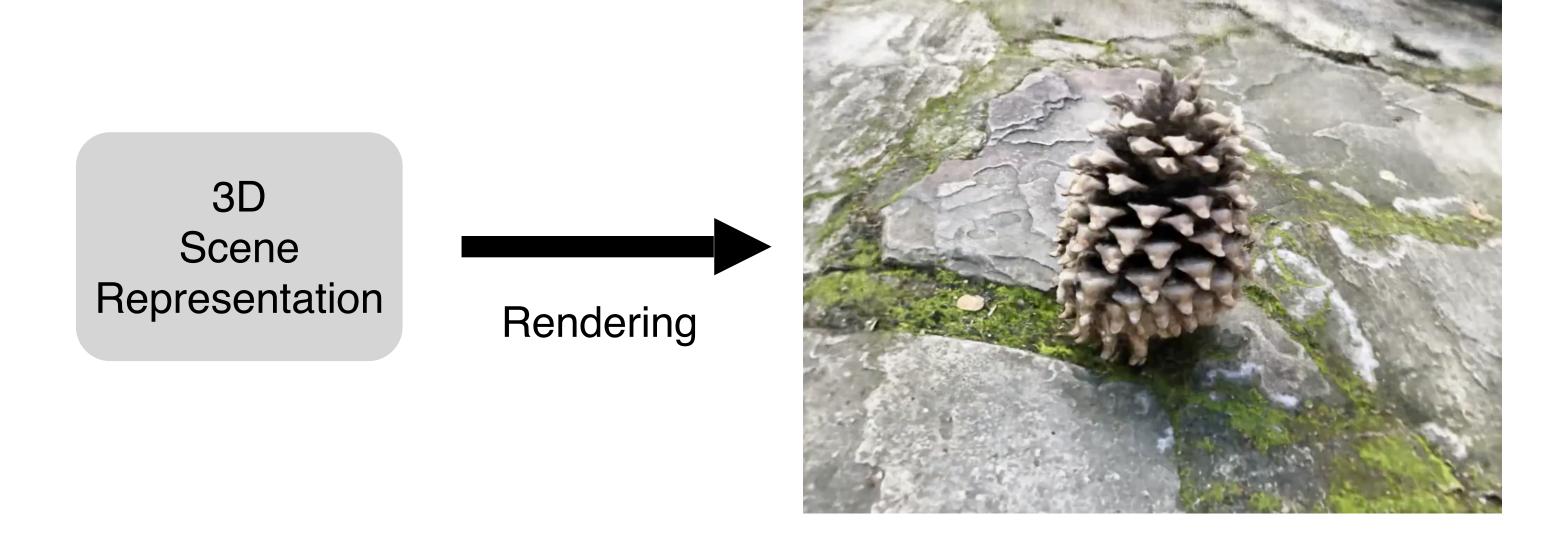
Lecture 22: Neural fields and learning-based 3D

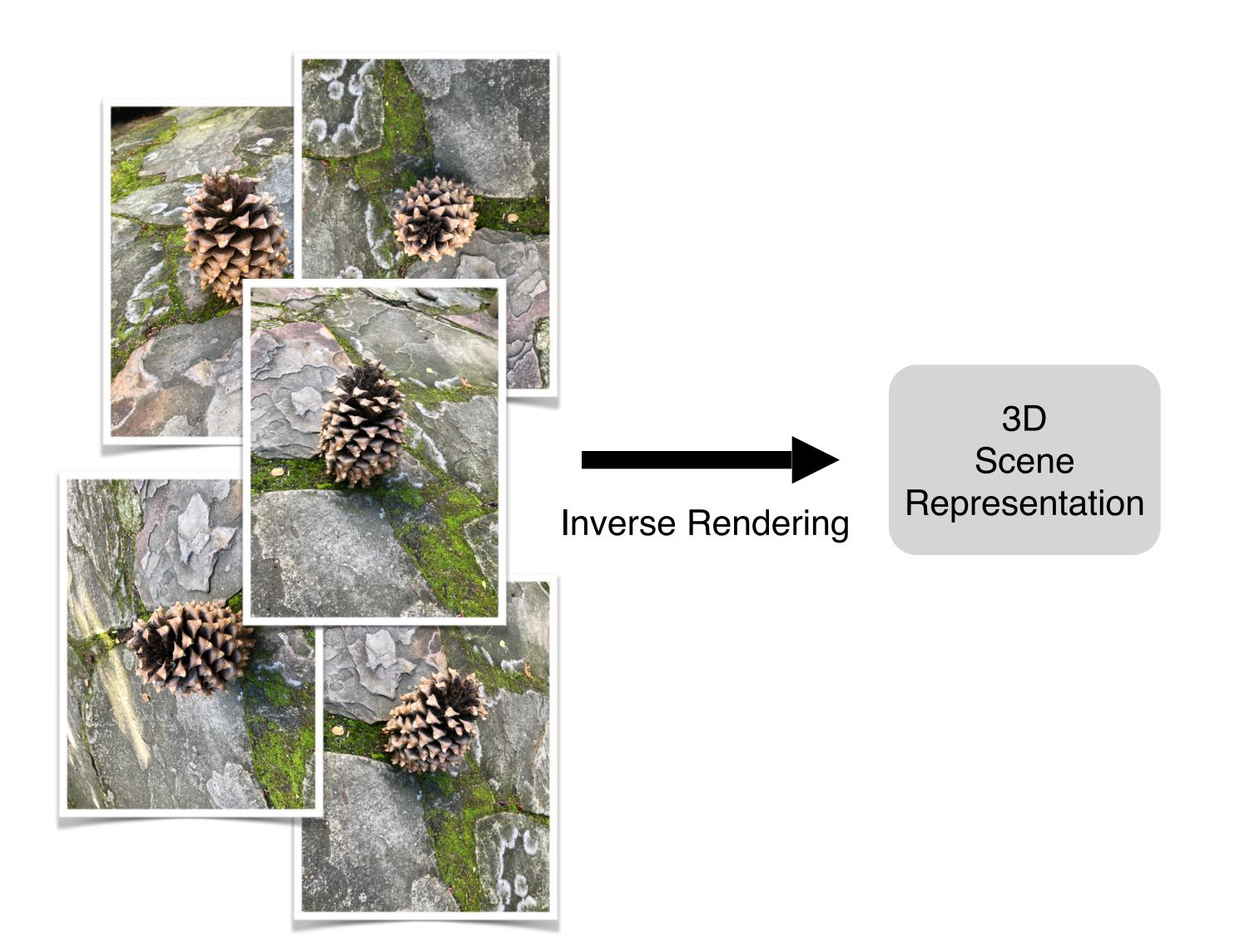
CS 5670: Introduction to Computer Vision



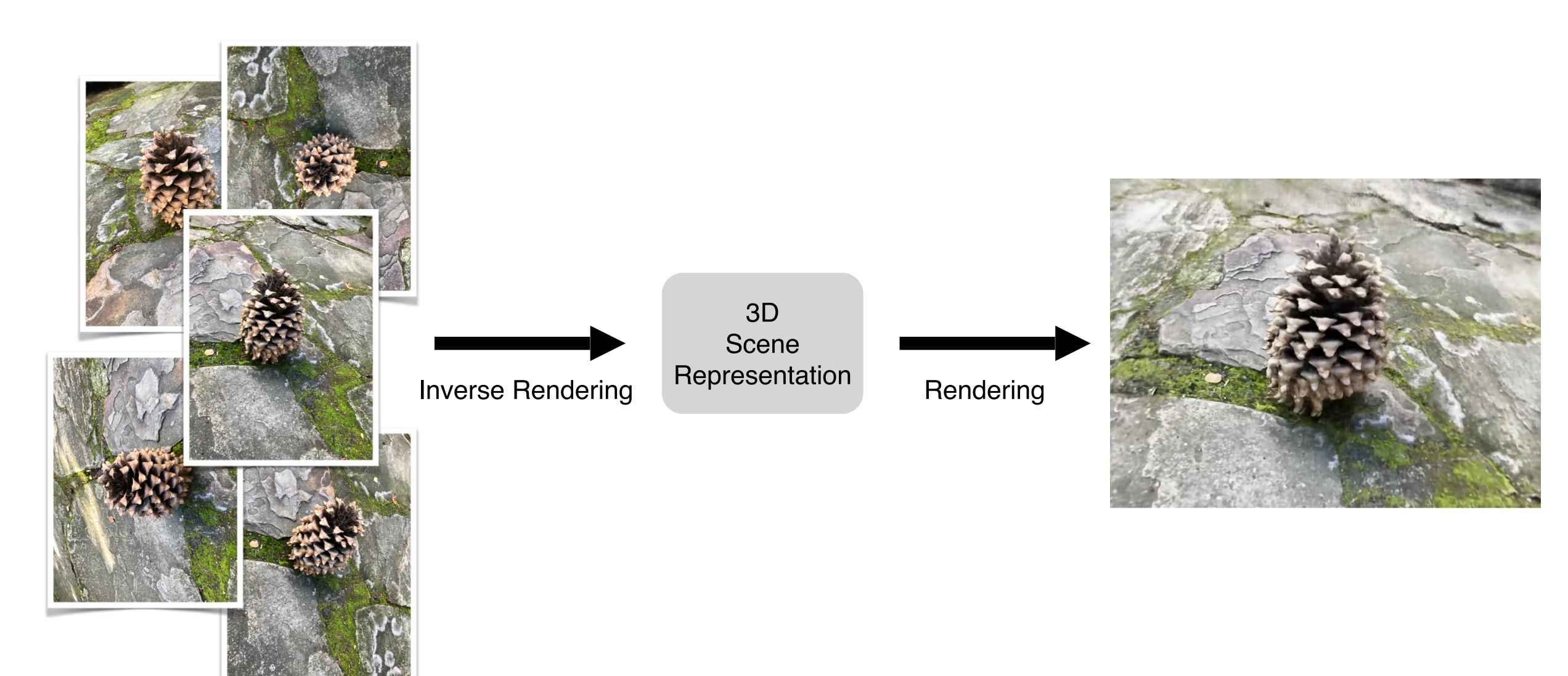
Rendering in computer graphics



Rendering in computer graphics



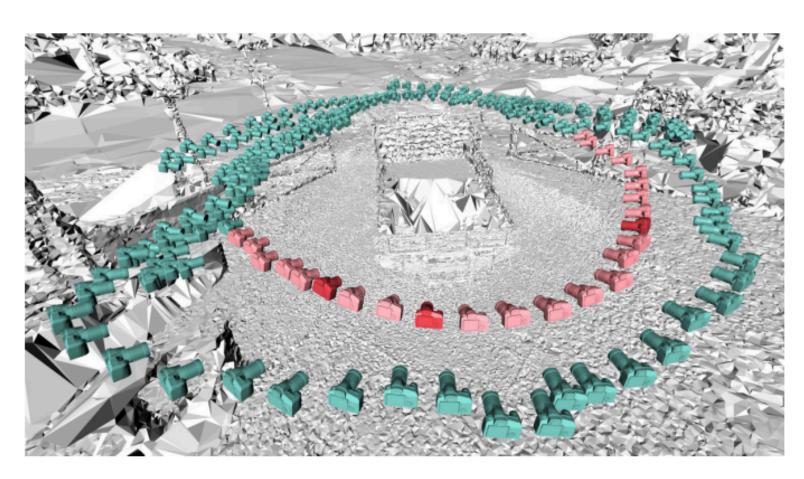
Rendering in computer graphics



Idea #1: Image-based rendering







Point cloud

Proxy geometry (a mesh)

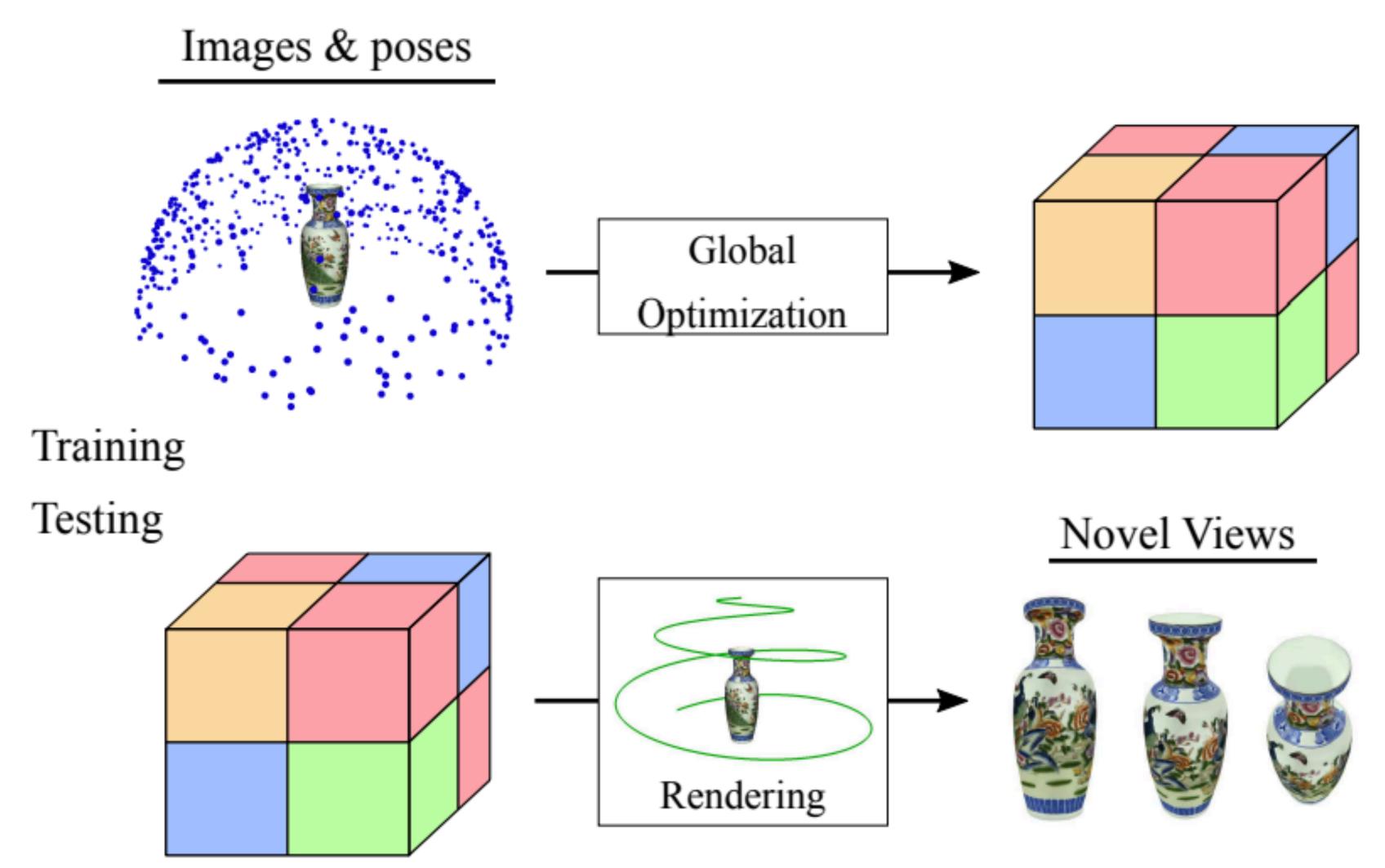
Reconstruction

To synthesize a new view, select colors from existing views using proxy geometry.

Idea #1: Image-based rendering

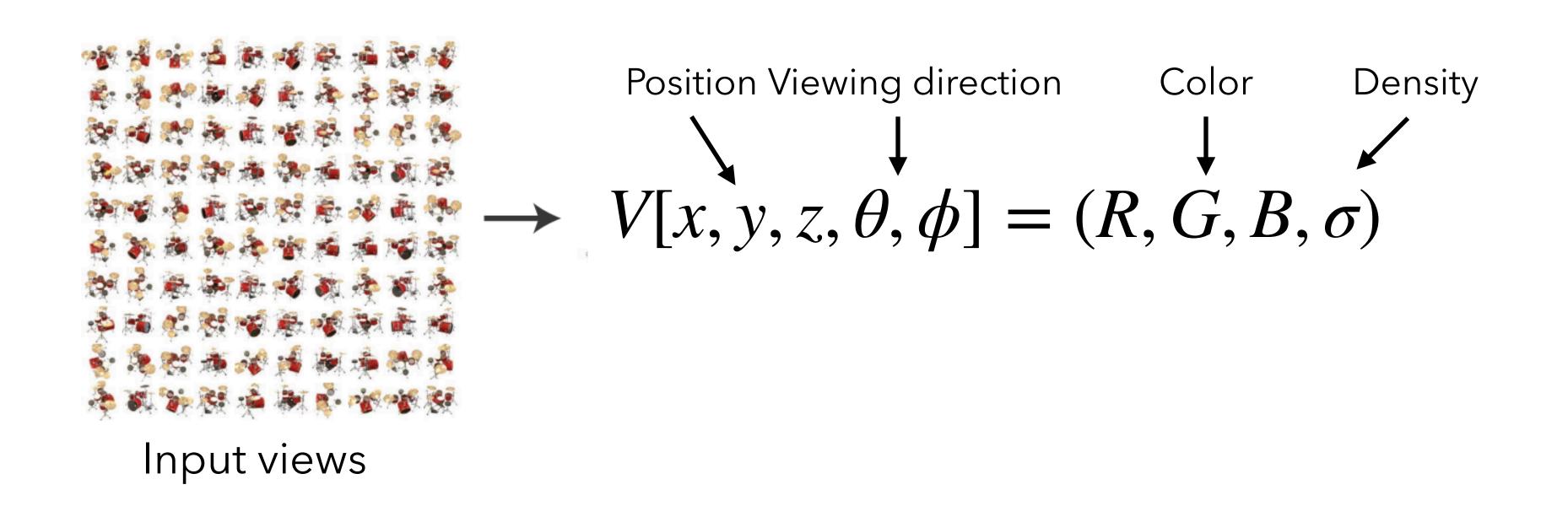


Idea #2: voxel representation

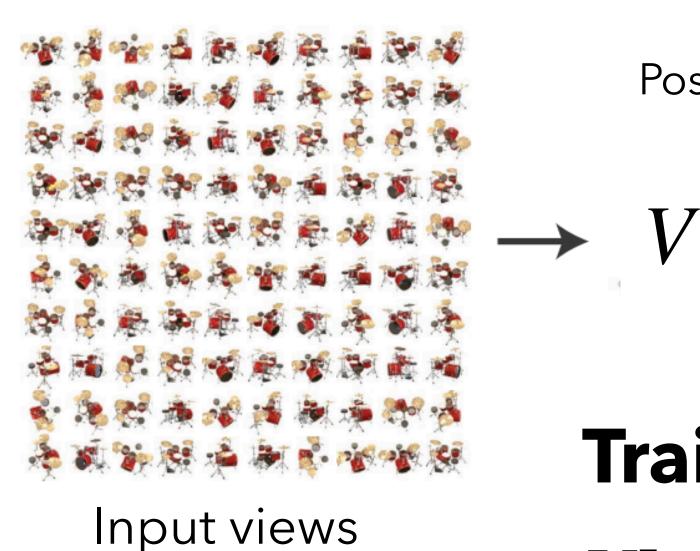


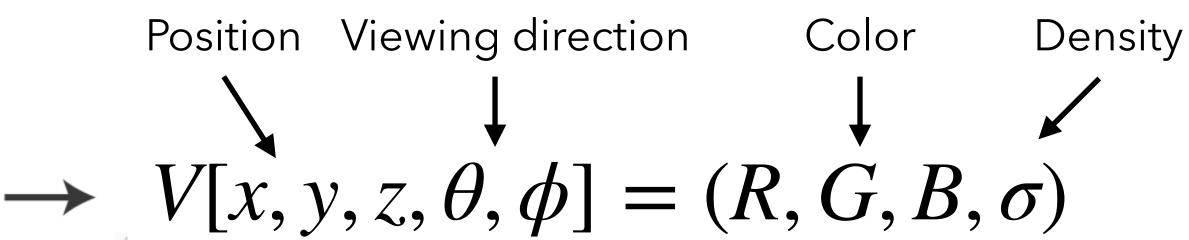
7

Idea #2: voxel representation

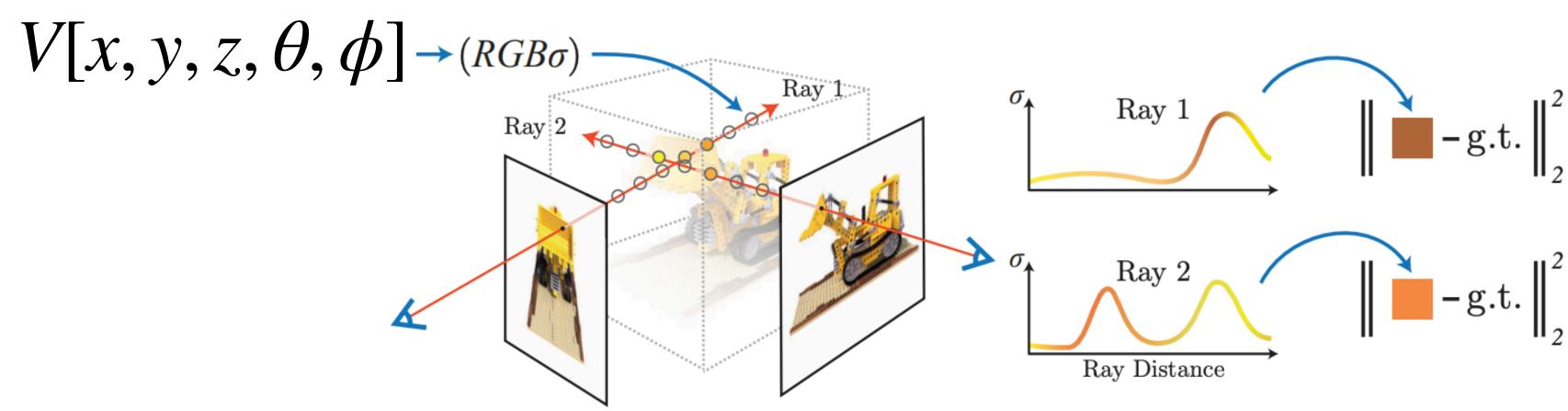


Idea #2: voxel representation



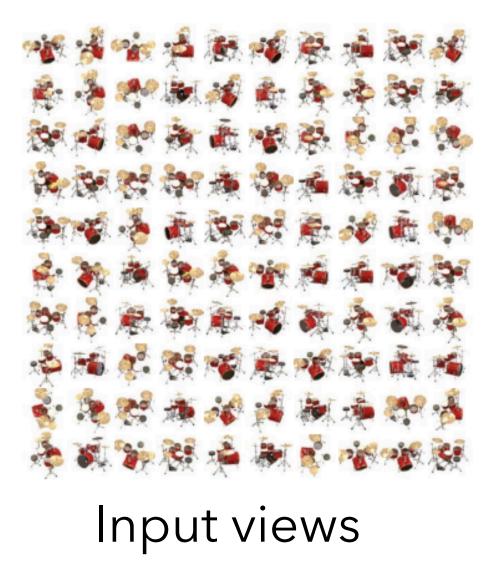


Training:



Problem: A huge table! $\mathcal{O}(D^3A^2)$

Idea #3: neural radiance field (NeRF)



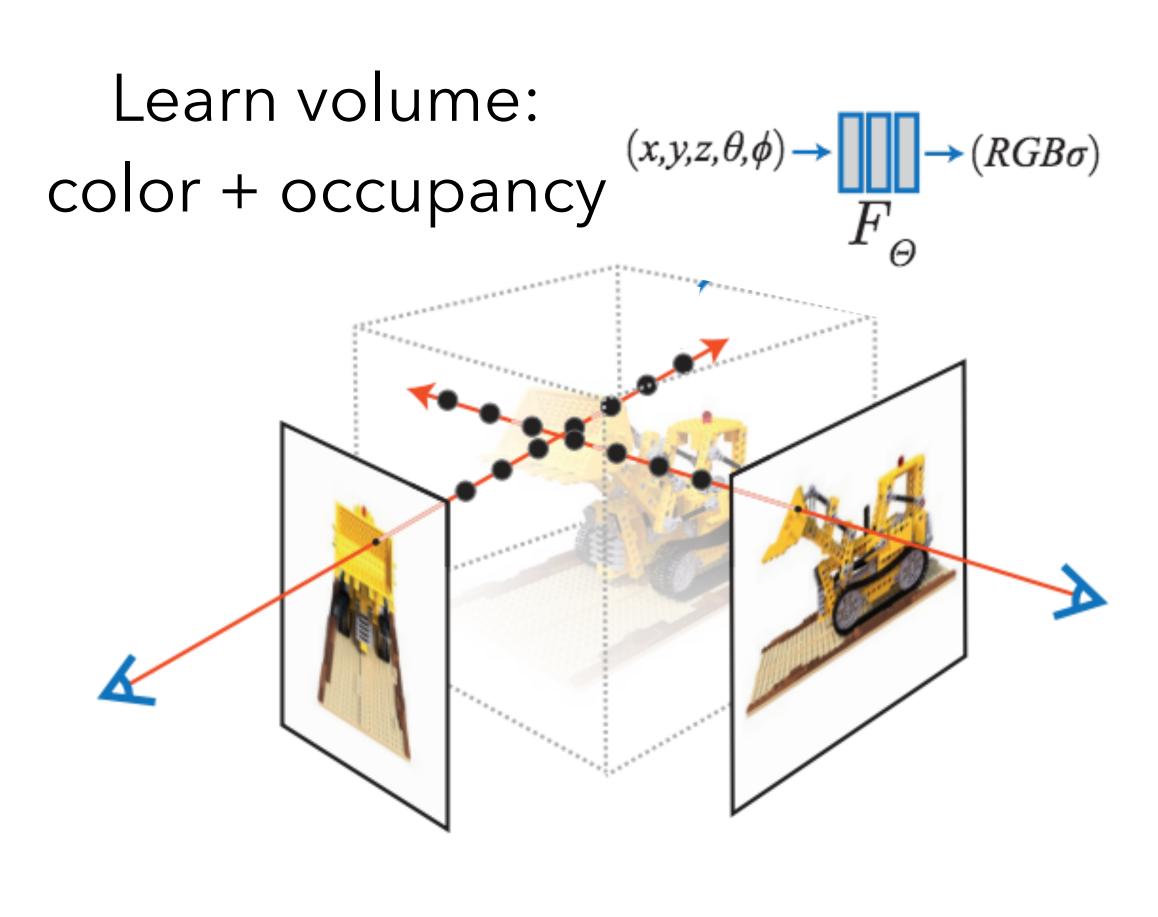
$$\rightarrow F_{\Theta}(x, y, z, \theta, \phi) = (R, G, B, \sigma)$$

- Represent using a neural radiance field.
- ullet Function that maps a (x, y, z, θ , ϕ) to a color and density.
- Typically parameterized as a multi-layer perceptron (MLP)
- ullet Goal: find parameters ullet for MLP that explain the images

Idea #3: neural radiance field (NeRF)

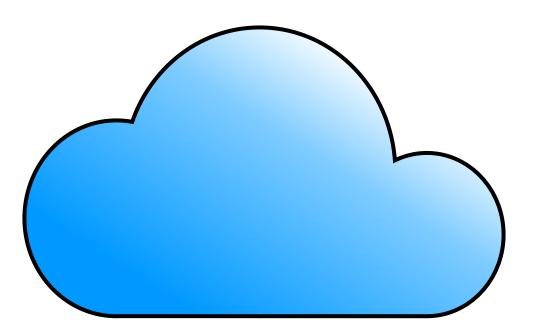


3D scene

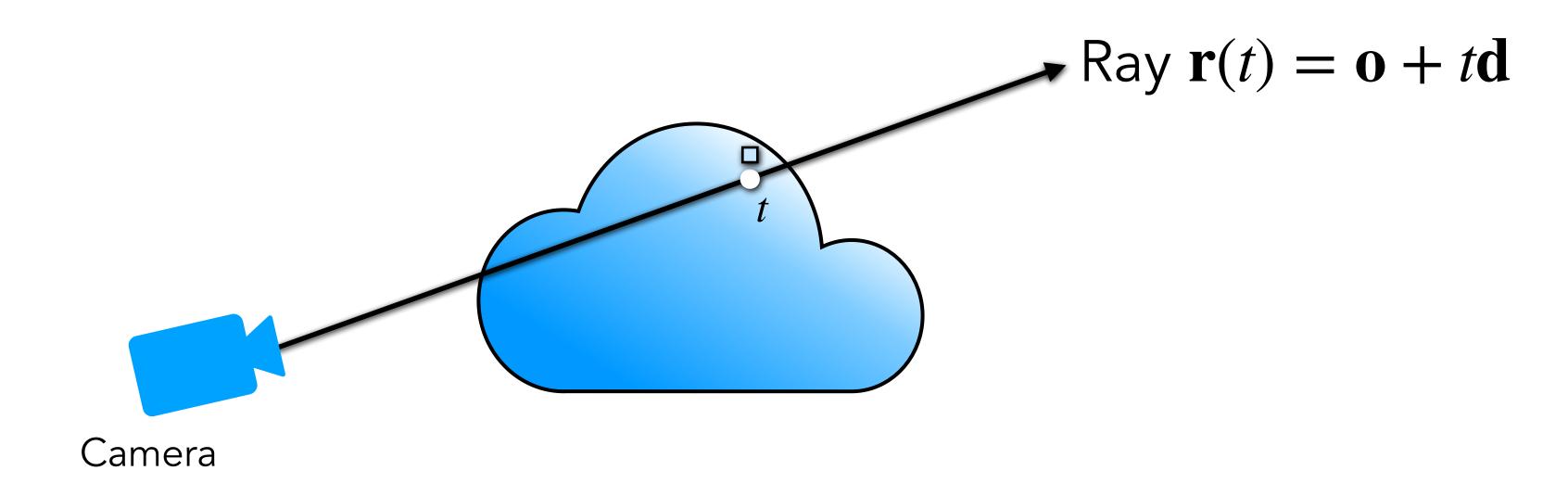


Viewpoints

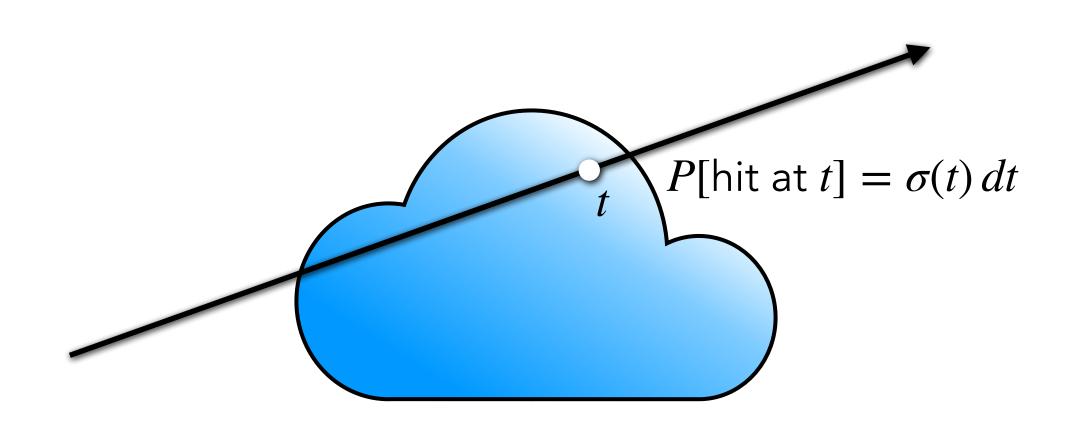
[Mildenhall*, Srinivasan*, Tanick*, et al., Neural radiance fields, 2020



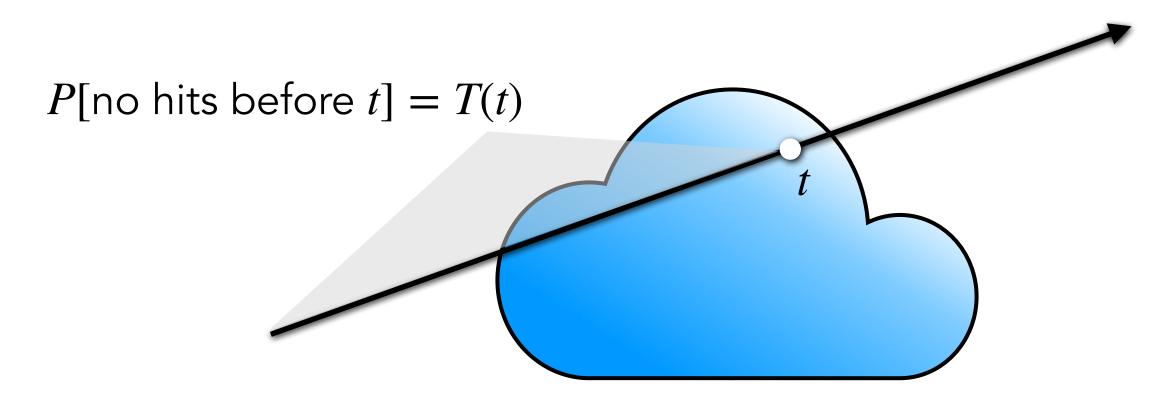
Scene is a cloud of tiny colored particles



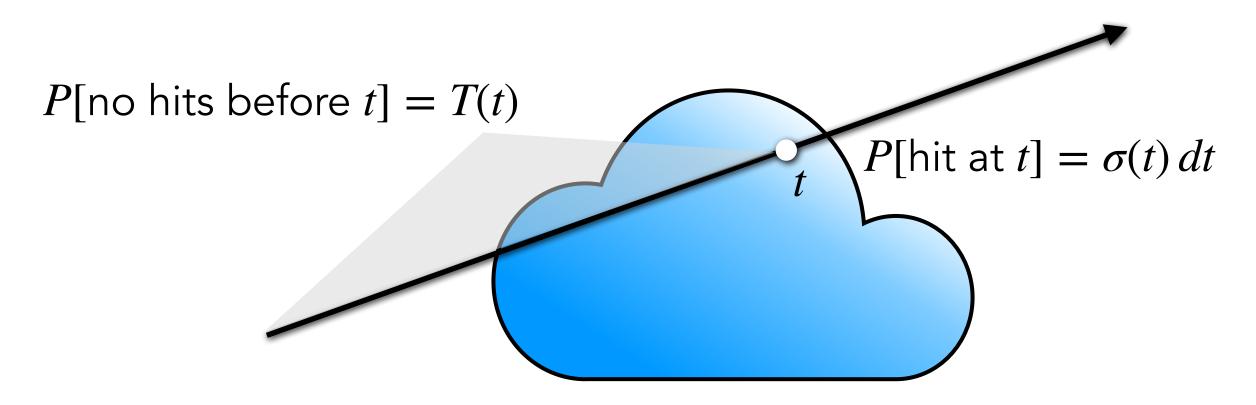
If a ray traveling through the scene hits a particle at t, we return its color $\mathbf{c}(t)$



Probability that ray stops in a small interval around t is $\sigma(t) dt$. σ is known as the **Volume Density.**

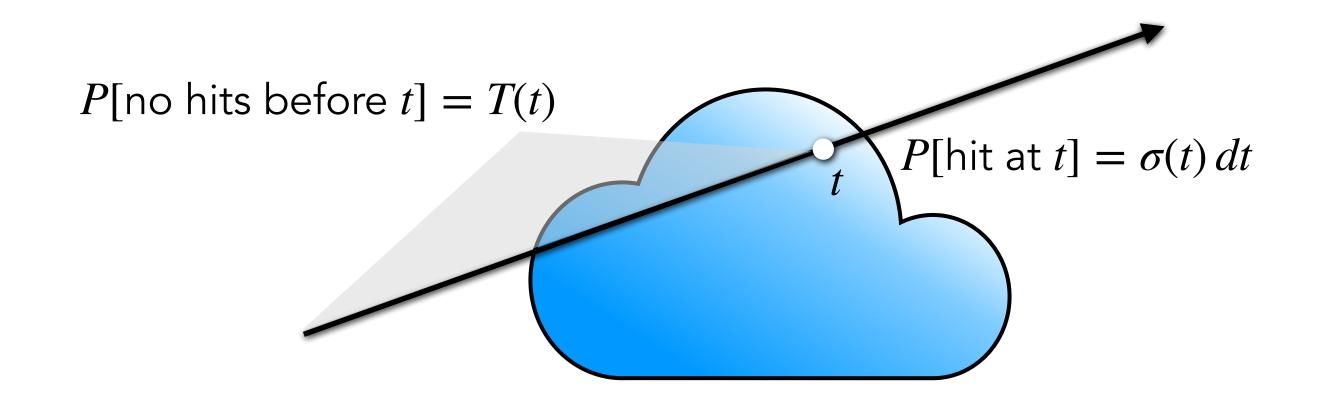


To determine if t is the *first* hit, need to know T(t): probability that the ray didn't hit any particles earlier. T(t) is called **Transmittance.**



 σ and T are related via

 $P[\text{no hit before } t + dt] = P[\text{no hit before } t] \times P[\text{no hit at } t]$



 $P[T(t+dt)] = P[T(t)] \times P[1-\sigma(t)dt]$

 σ and T are related via

Can show: no hits before t is equal to integral over density up until t.

$$T(t) = \exp\left(-\int_0^t \sigma(a)da\right)$$

1 - T(t) can be seen as cumulative distribution function of probability that ray hits something before reaching t.

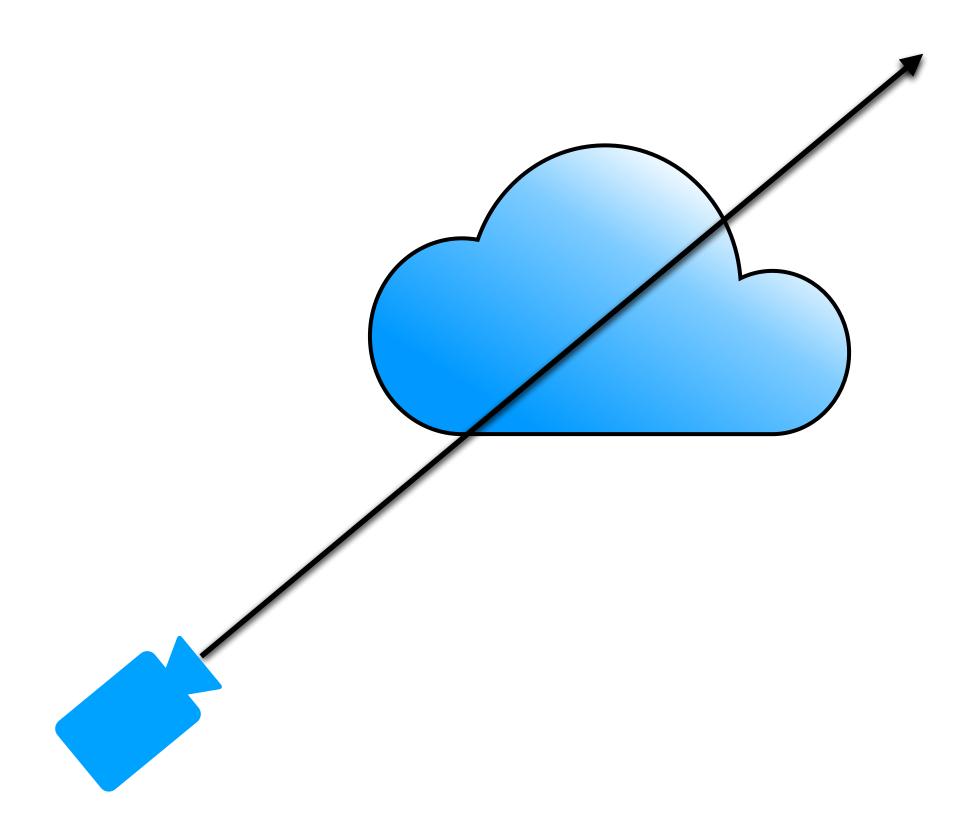
Then $T'(t) = T(t)\sigma(t)$ is probability that ray stops exactly at t.

Then $T'(t) = T(t)\sigma(t)$ is probability that ray stops exactly at t.

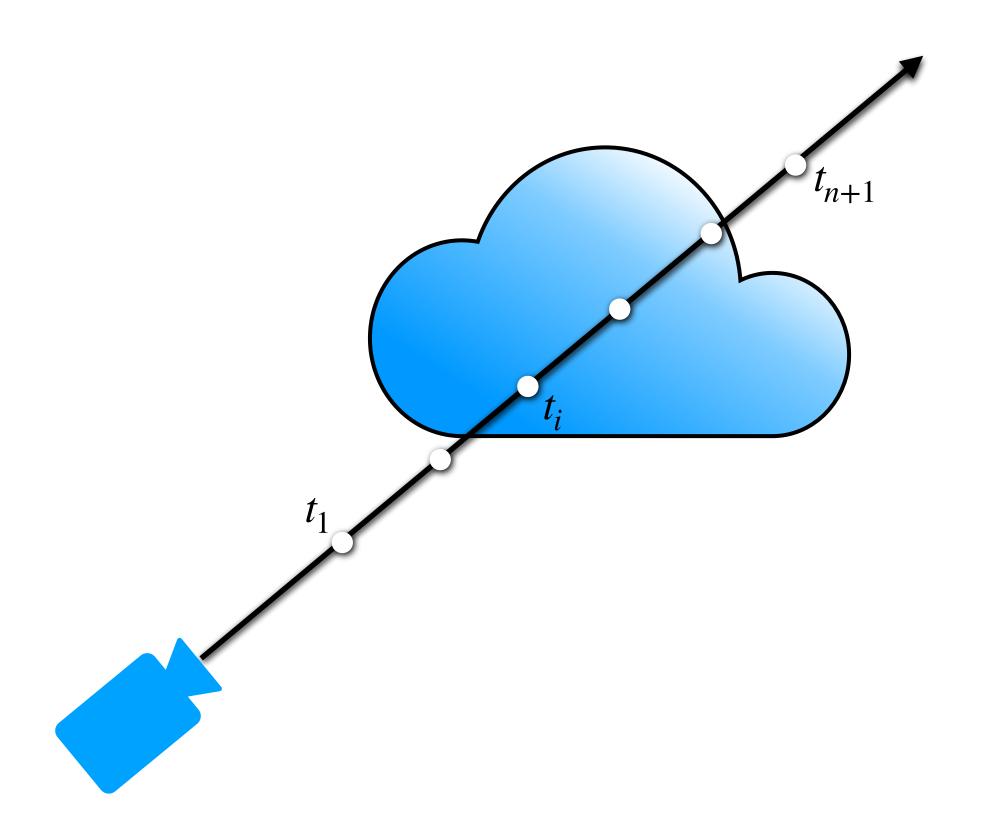
So the expected color returned by the ray will be

$$\int_{t_0}^{t_1} T(t)\sigma(t)\mathbf{c}(t) dt$$

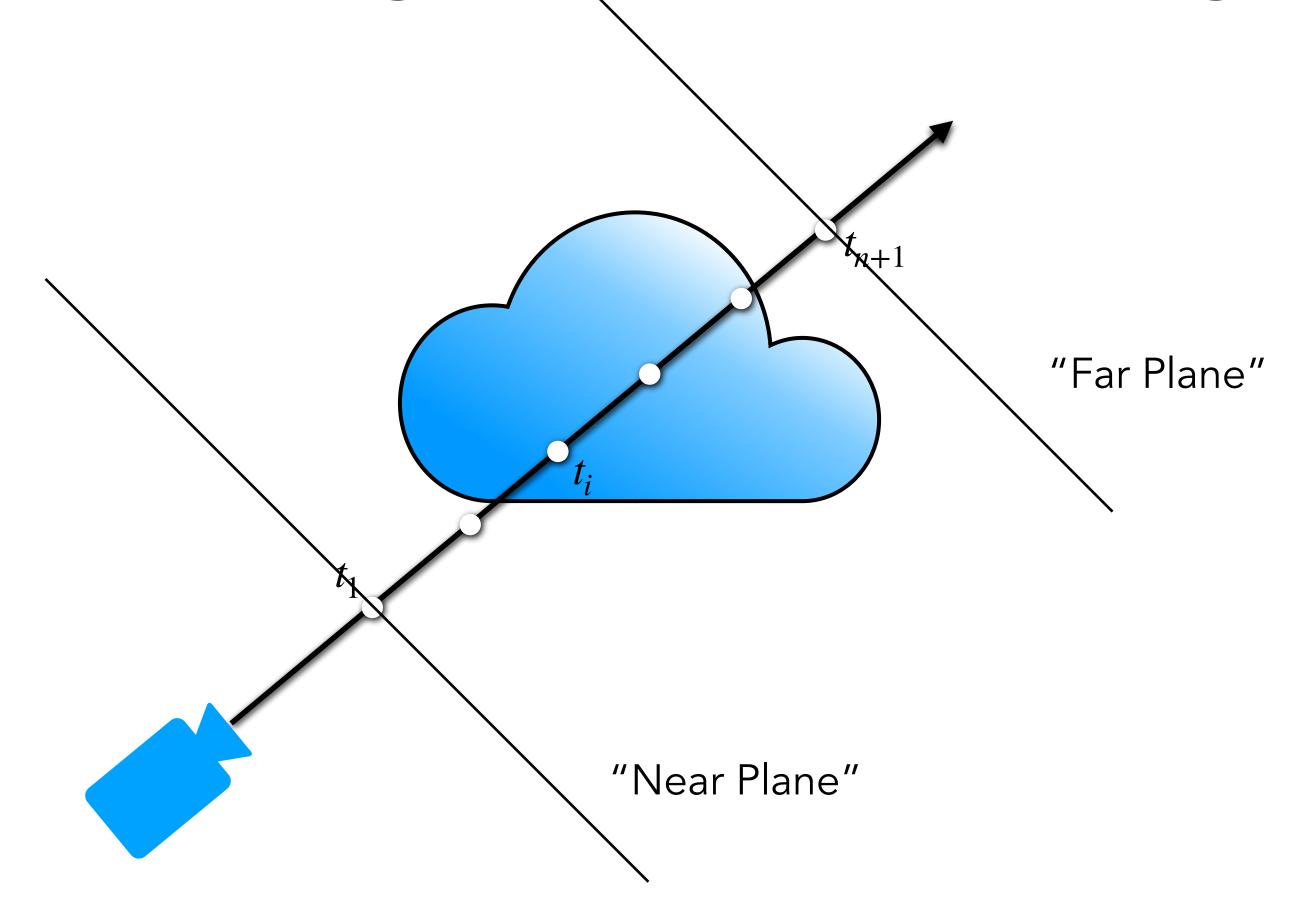
Note the nested integral!



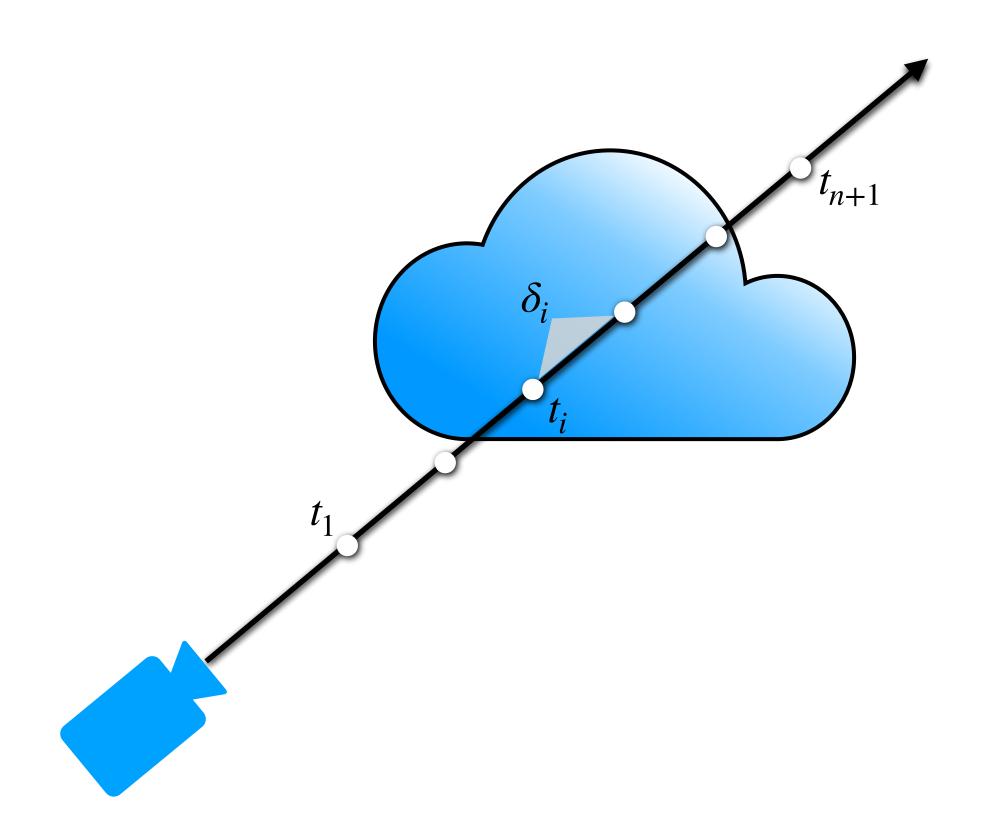
We use quadrature to approximate the nested integral,



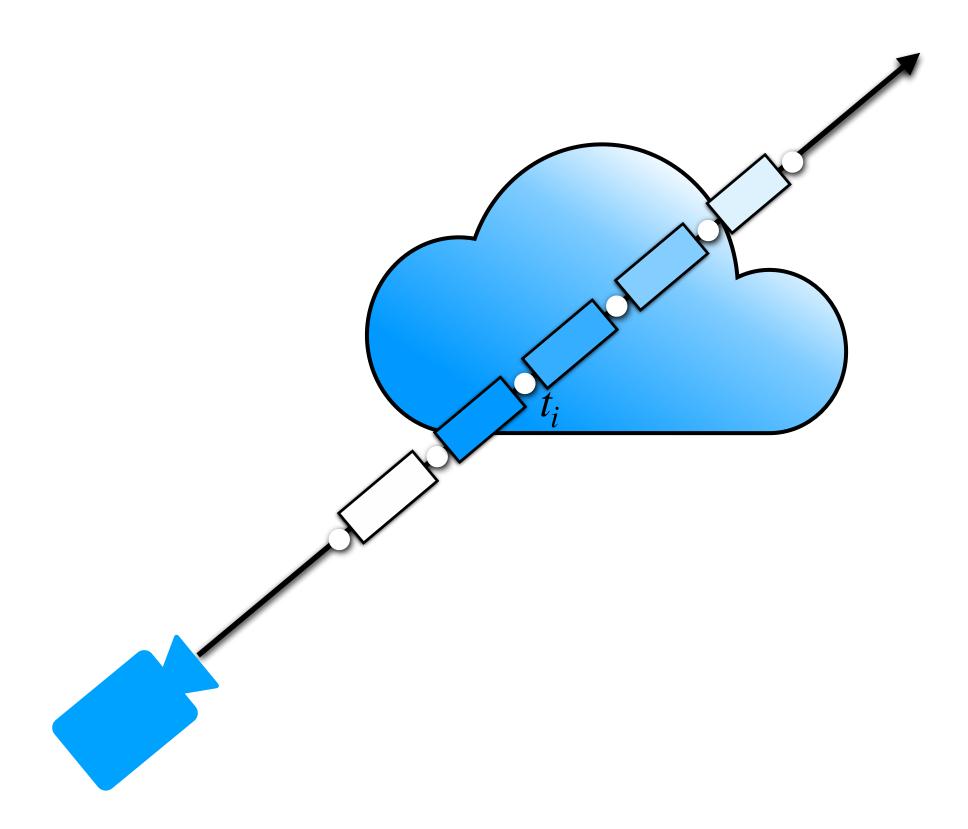
We use quadrature to approximate the nested integral, splitting the ray up into n segments with endpoints $\{t_1, t_2, ..., t_{n+1}\}$



We use quadrature to approximate the nested integral, splitting the ray up into n segments with endpoints $\{t_1, t_2, ..., t_{n+1}\}$



We use quadrature to approximate the nested integral, splitting the ray up into n segments with endpoints $\{t_1, t_2, ..., t_{n+1}\}$ with lengths $\delta_i = t_{i+1} - t_i$



We assume volume density and color are roughly constant within each interval

Summary: volume rendering integral estimate

Rendering model for ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$:

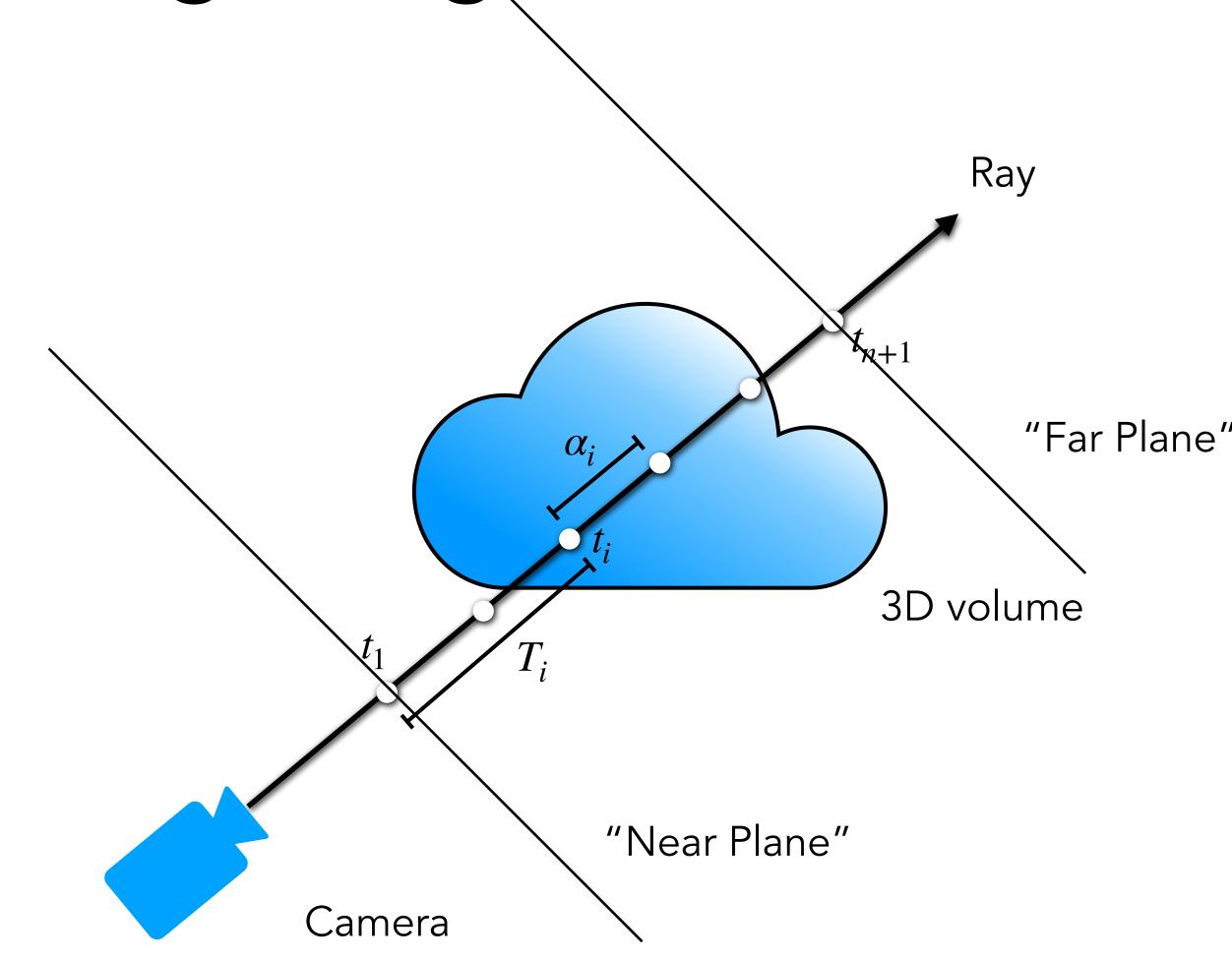
$$\mathbf{c} pprox \sum_{i=1}^{n} T_i \alpha_i \mathbf{c}_i$$
 \mathbf{c}_i
 \mathbf{c}_i

How much light is blocked earlier along ray:

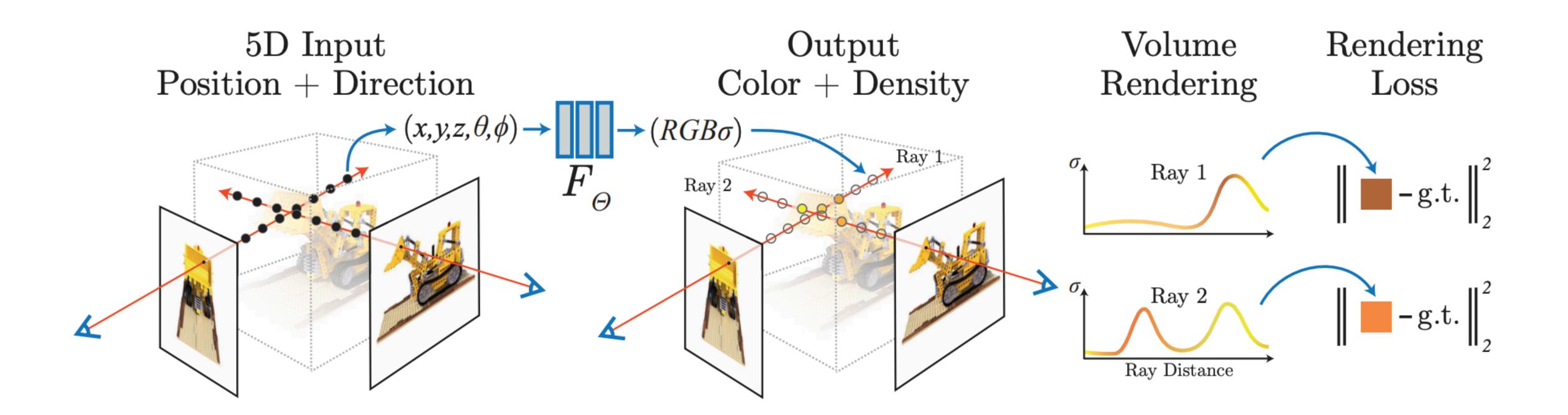
$$T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$$

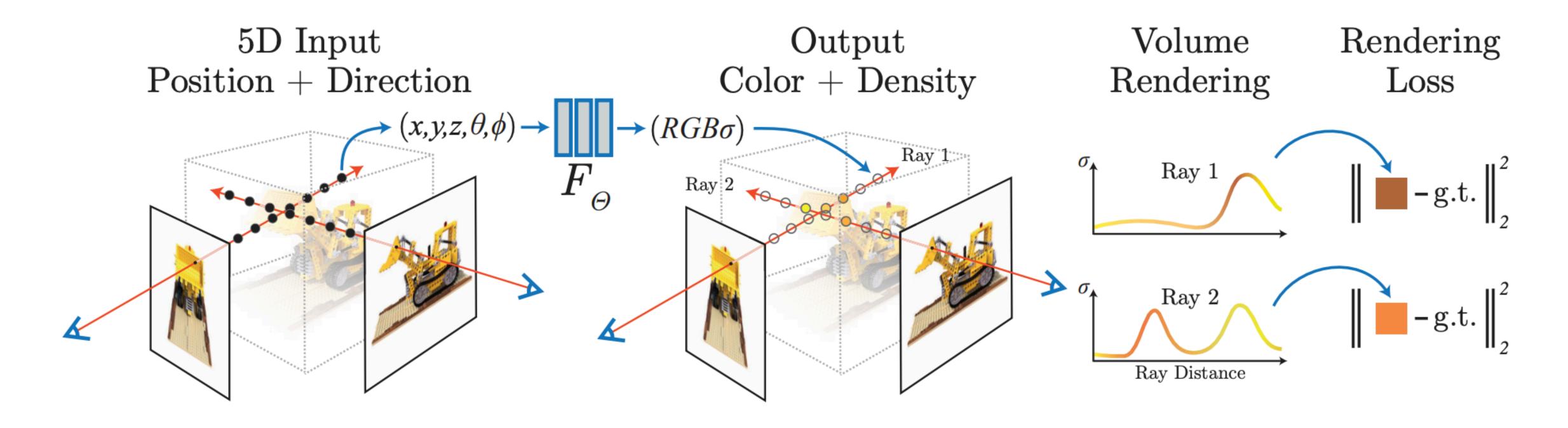
How much light is contributed by ray segment *i*:

$$\alpha_i = 1 - \exp(-\sigma_i \delta_i)$$



Learning a NeRF





$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))\mathbf{c}(\mathbf{r}(t), \mathbf{d})dt, \text{ where } T(t) = \exp\left(-\int_{t_n}^{t} \sigma(\mathbf{r}(s))ds\right)$$

Ray: $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ For color c and density σ .

Color for ray **r**

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))\mathbf{c}(\mathbf{r}(t),\mathbf{d})dt, \text{ where } T(t) = \exp\left(-\int_{t_n}^{t} \sigma(\mathbf{r}(s))ds\right)$$

Ray:
$$\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$$
 For color c and density σ .

A point distance t along \mathbf{r} , centered at

Color for ray **r**

Weight

Color at 3D point $\mathbf{r}(t)$

& direction d

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt, \text{ where } T(t) = \exp\left(-\int_{t_n}^{t} \sigma(\mathbf{r}(s)) ds\right)$$

Ray: $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$

For color c and density σ .

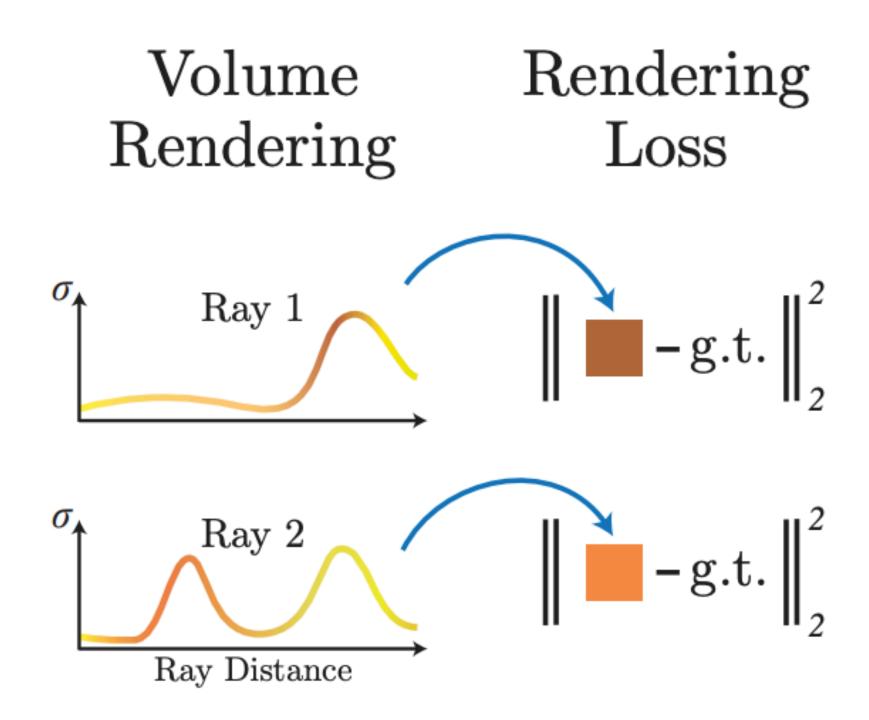
Density at point $\mathbf{r}(t)$

been absorbed $C(\mathbf{r}) = \int_{t}^{t_f} T(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt$, where $T(t) = \exp\left(-\int_{t}^{t} \sigma(\mathbf{r}(s)) ds\right)$

For color c and density σ . Ray: $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$

Probability that ray hasn

Loss function



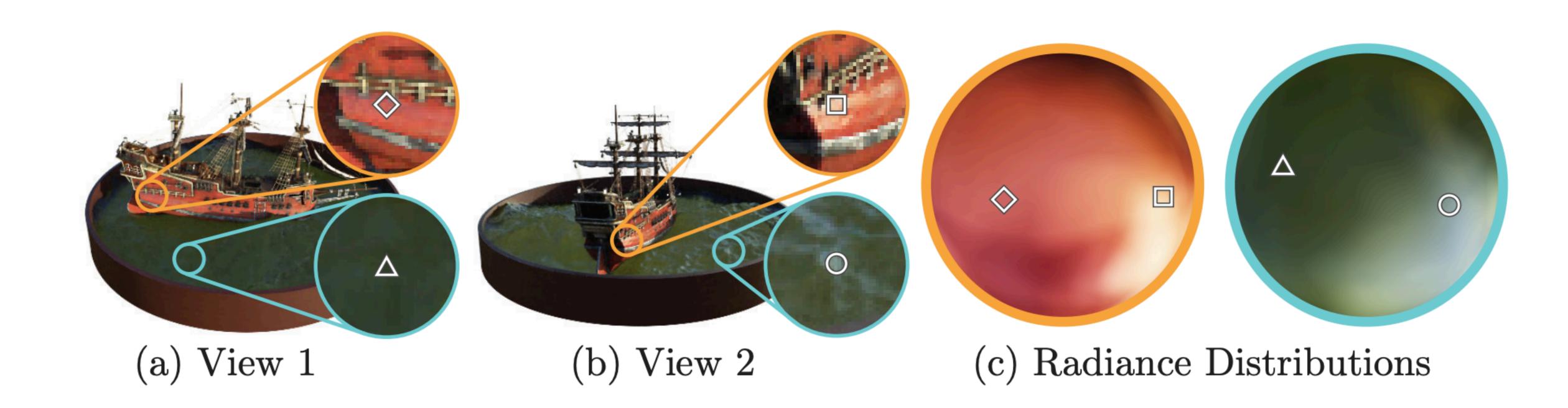
$$\mathcal{L} = \sum_{\mathbf{r} \in \mathcal{R}} \|C(\mathbf{r}) - C_{gt}(\mathbf{r})\|_{2}^{2}$$

Minimize difference between predicted and observed colors.

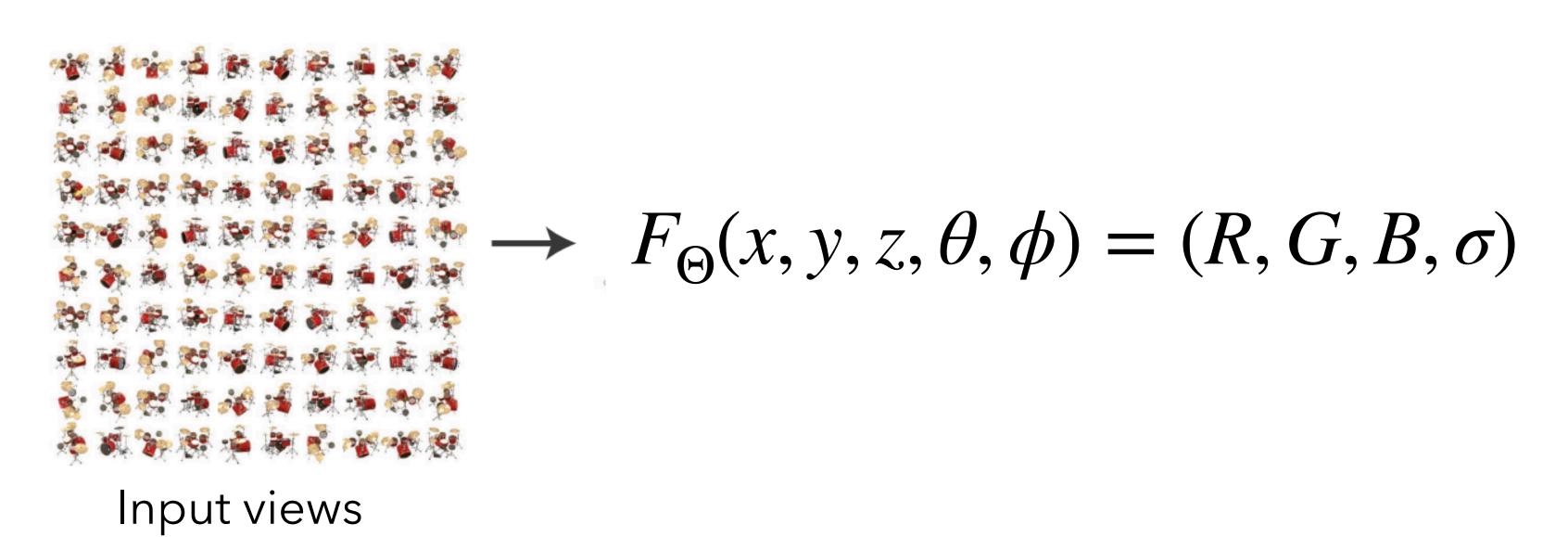
In practice: coarse-to-fine and other tricks.

Implementation details

Why is it good to be view-dependent?



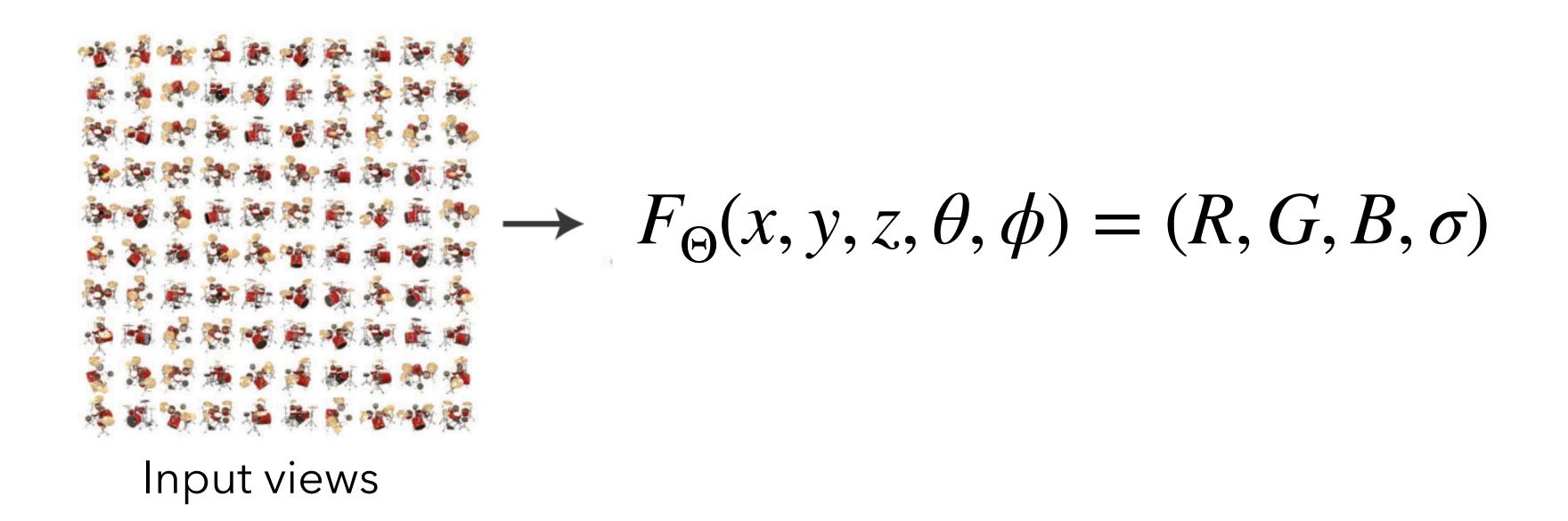
Representing the inputs



- In theory, could just plug in 4 inputs x, y, z, θ, ϕ
- However, this leads to blurry results.
- Neural nets show a bias toward low frequency functions [Tancik et al., 2020]



Fourier features

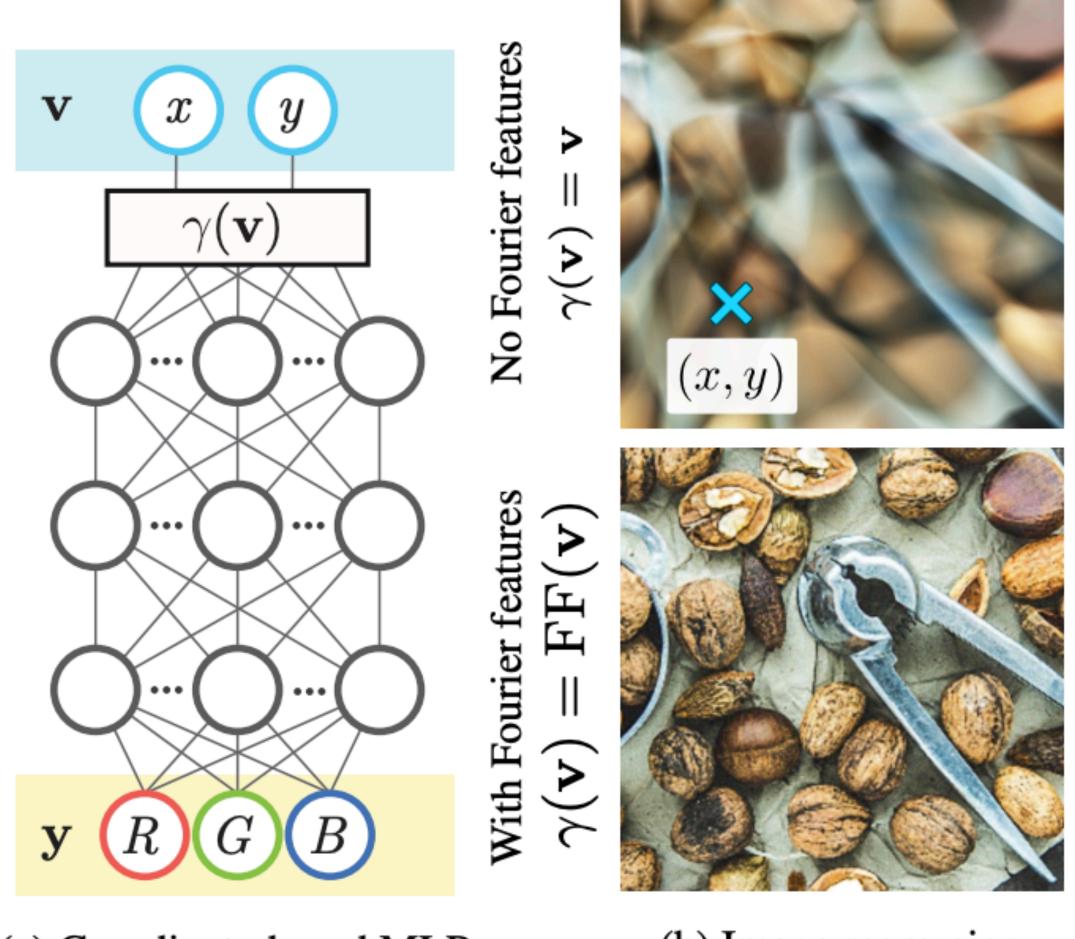


ullet Use a **positional encoding**. Given a scalar p, compute:

$$\gamma(p) = (\sin(2^0 \pi p), \cos(2^0 \pi p), \dots, \sin(2^{L-1} \pi p), \cos(2^{L-1} \pi p))$$

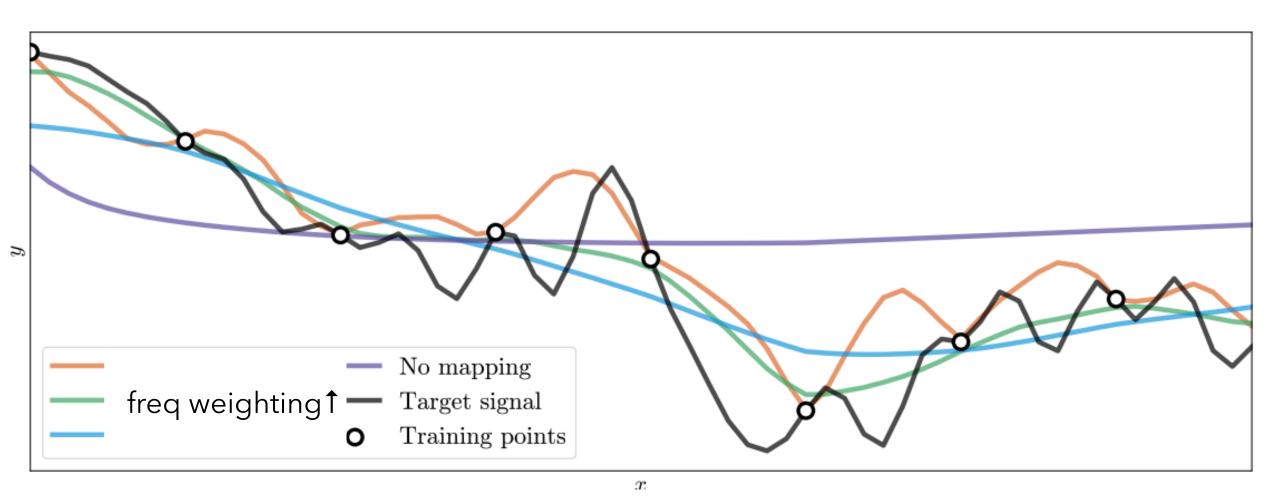
 Plug in the coordinate to sinusoids at different frequencies (e.g. L = 10).

Why does this happen?



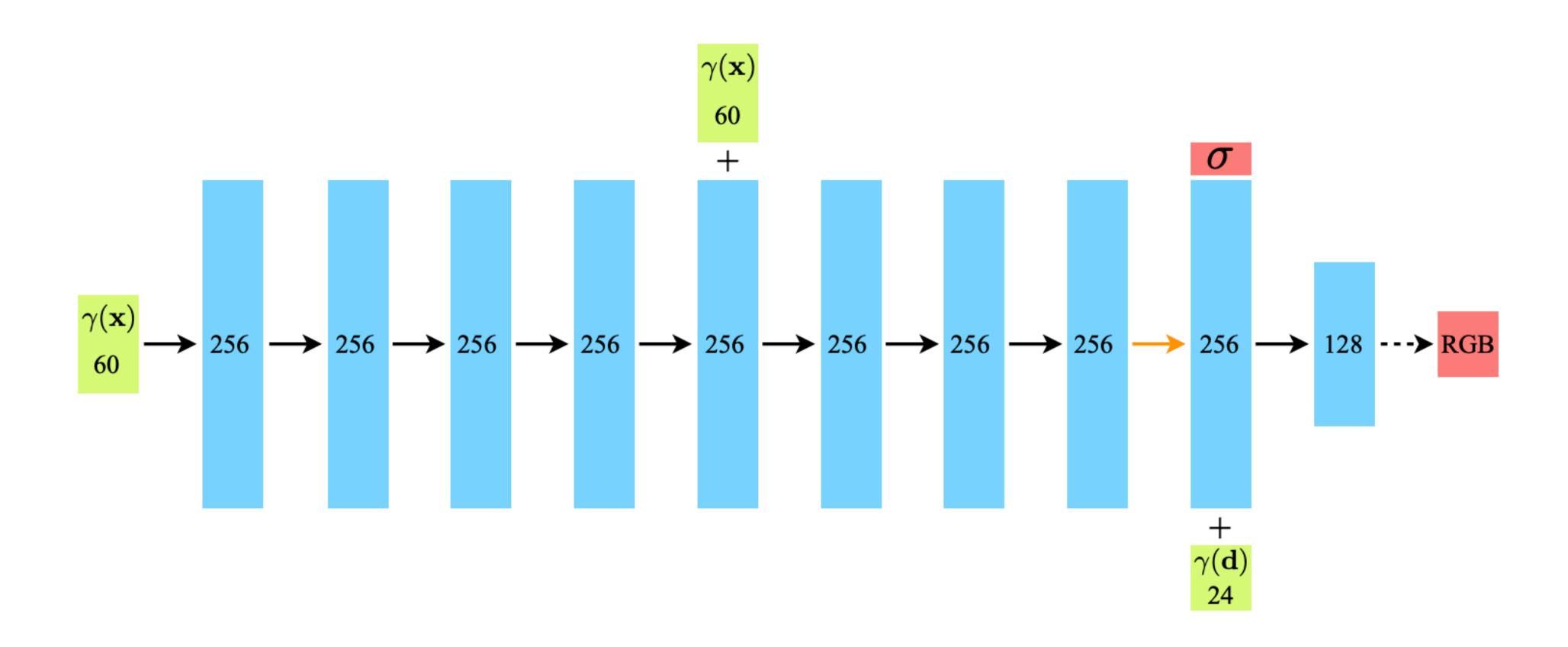
(a) Coordinate-based MLP (b) Image regression $(x,y) \to RGB$

- Neural nets have trouble learning high frequency functions
- This mapping explicitly represents
 different frequencies (forces net to pay
 more attention high frequencies)



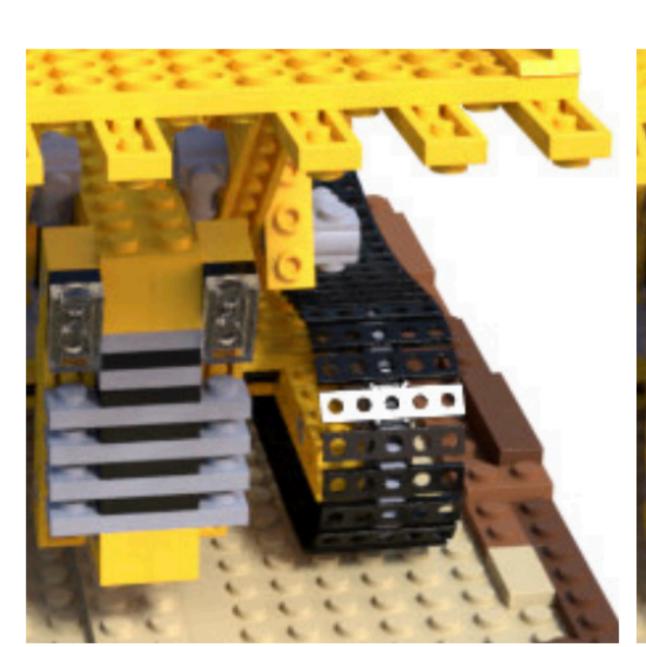
See [Tancik et al., "Fourier Features Let Networks Learn High Frequency Functions in Low Dimensional Domains", 2020]

MLP architecture

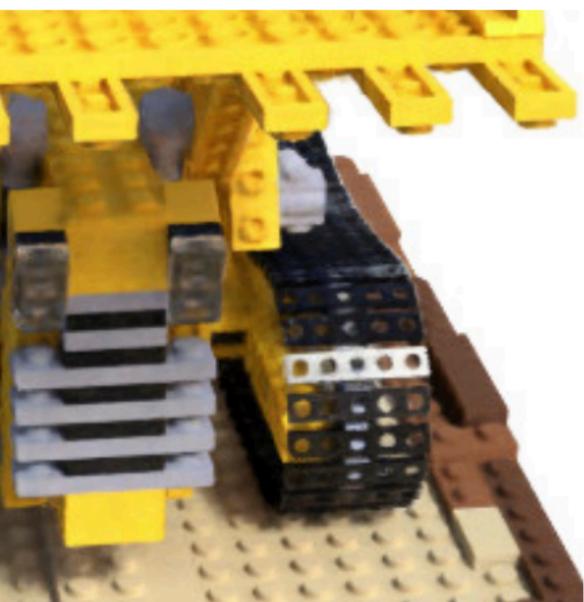


$$\gamma(p) = (\sin(2^0\pi p), \cos(2^0\pi p), \dots, \sin(2^{L-1}\pi p), \cos(2^{L-1}\pi p))$$

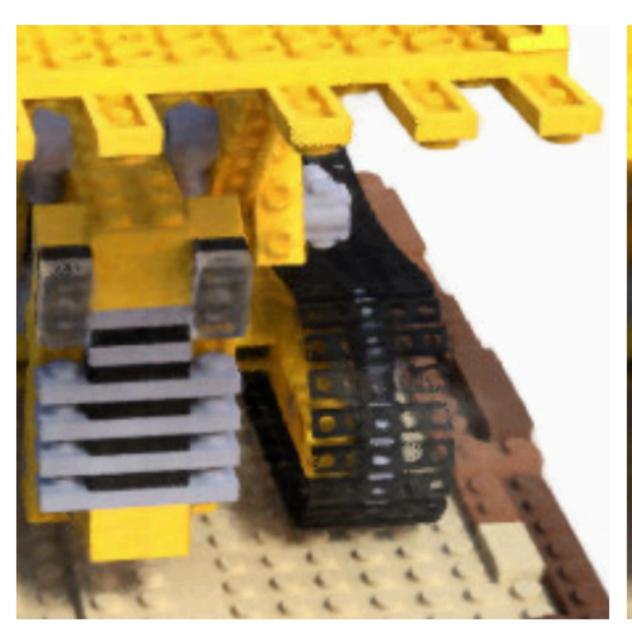
Results for a novel viewpoint



Ground Truth



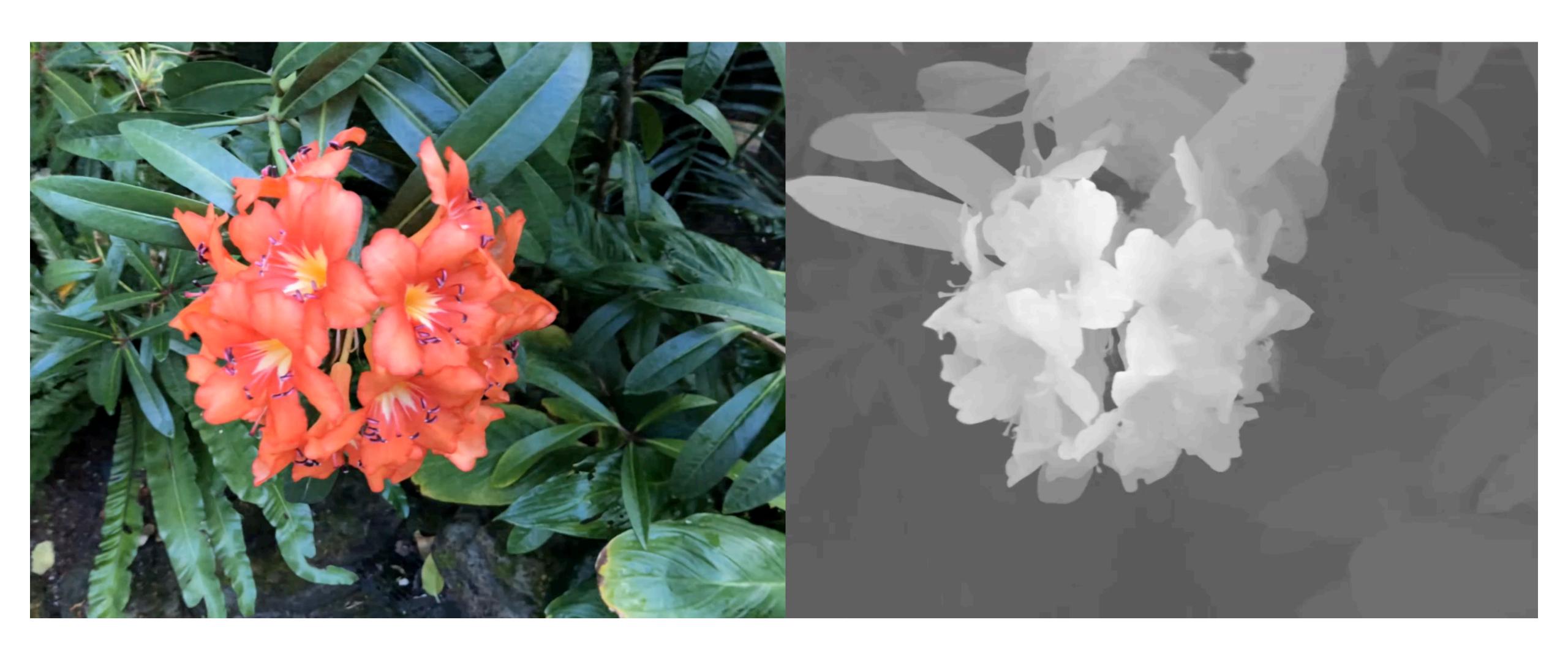
Complete Model





No View Dependence No Positional Encoding

Results



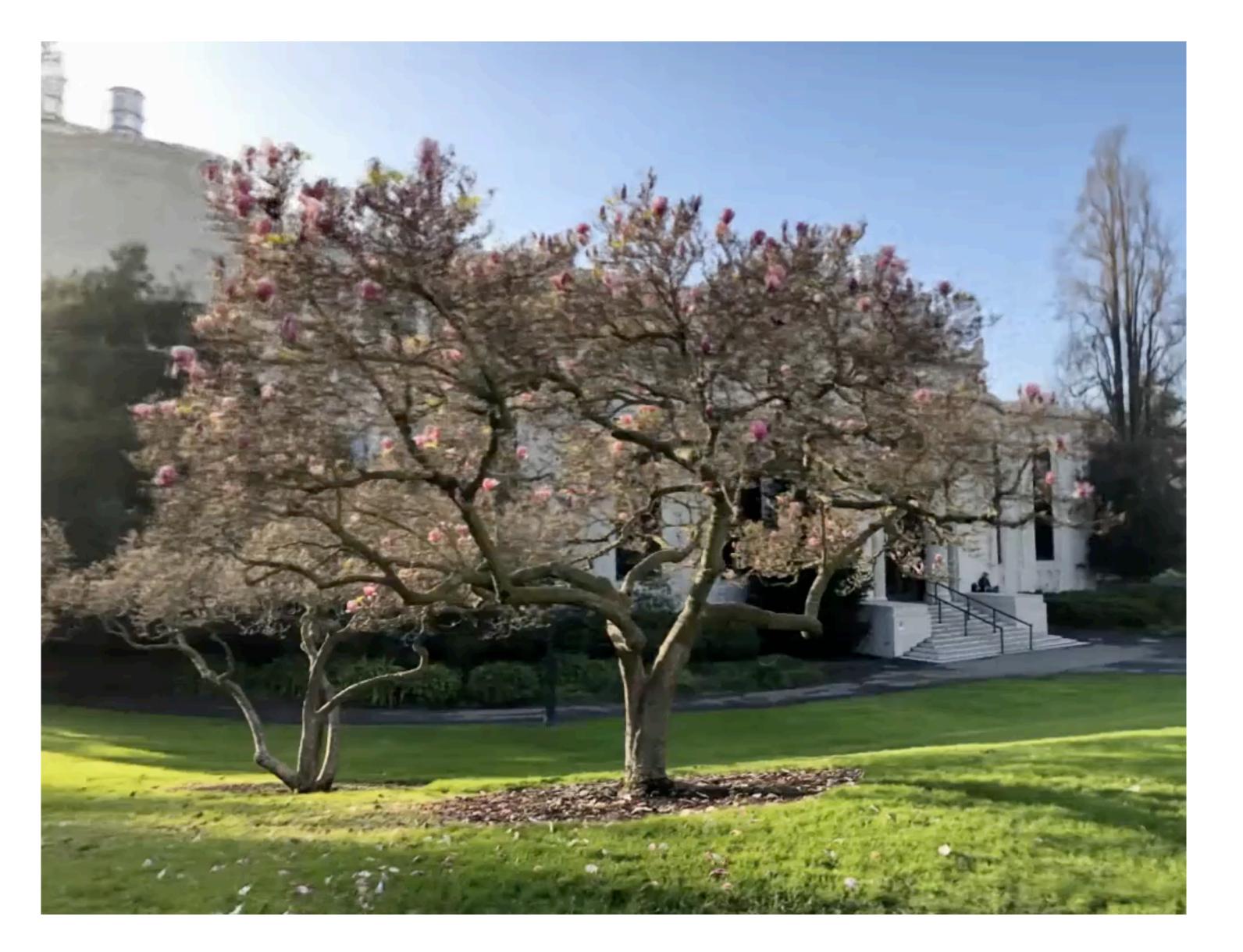
[Mildenhall*, Srinivasan*, Tanick*, et al. 2020

Results



[Mildenhall*, Srinivasan 40 Tanick*, et al. 2020

Results



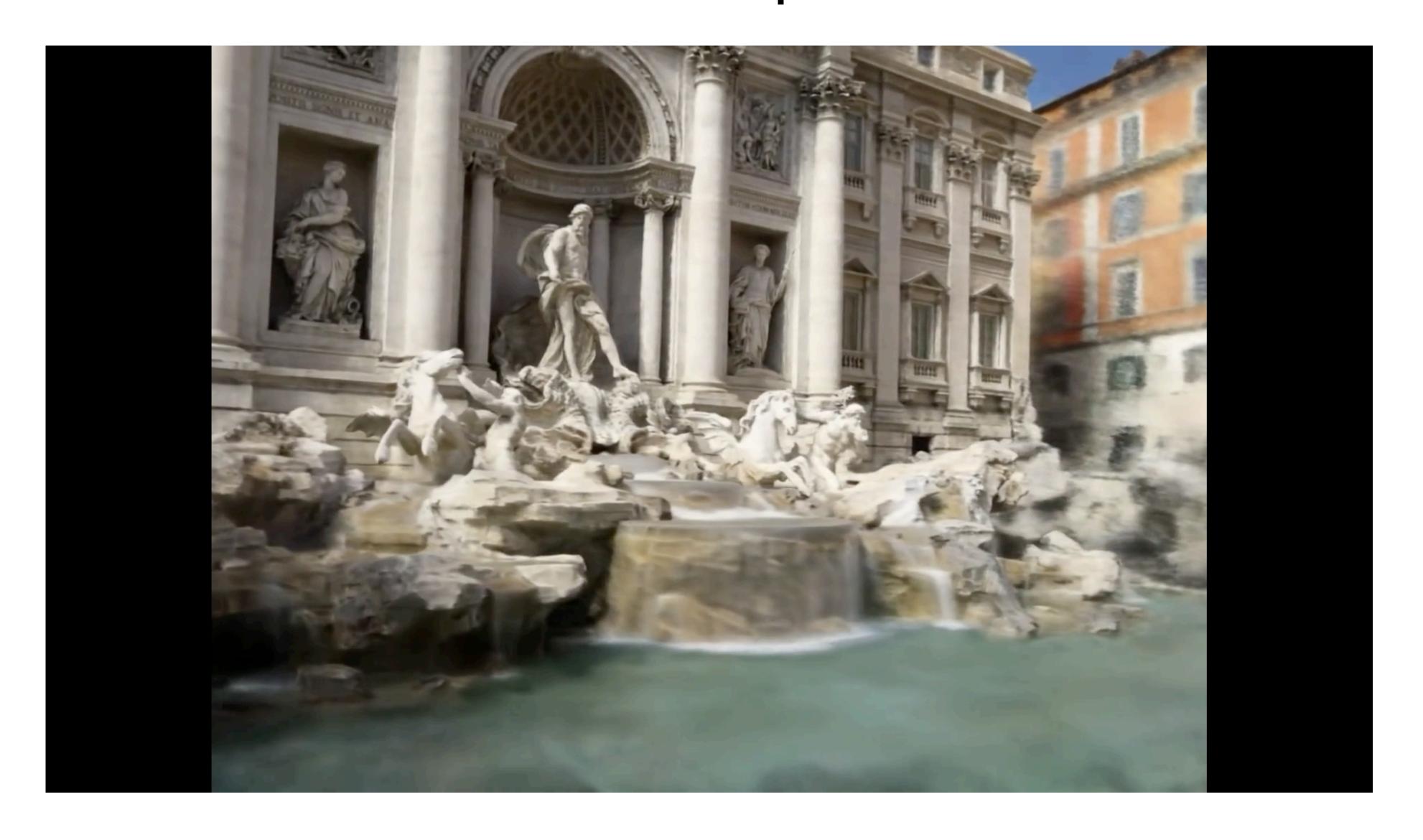
[Mildenhall*, Srinivasan*, Tanick*, et al. 2020

NeRF with other tricks



[Barron et al., "Zip-NeRF: Anti-Aliased Grid-Based Neural Radiance Fields. 2023]

Extension: internet photo collections



Extension: internet photo collections





Matthew Tancik1*

Vincent Casser²

Xinchen Yan²

Sabeek Pradhan²

Ben Mildenhall³

Pratul Srinivasan³

Jonathan T. Barron³ Henrik Kretzschmar²

¹UC Berkeley

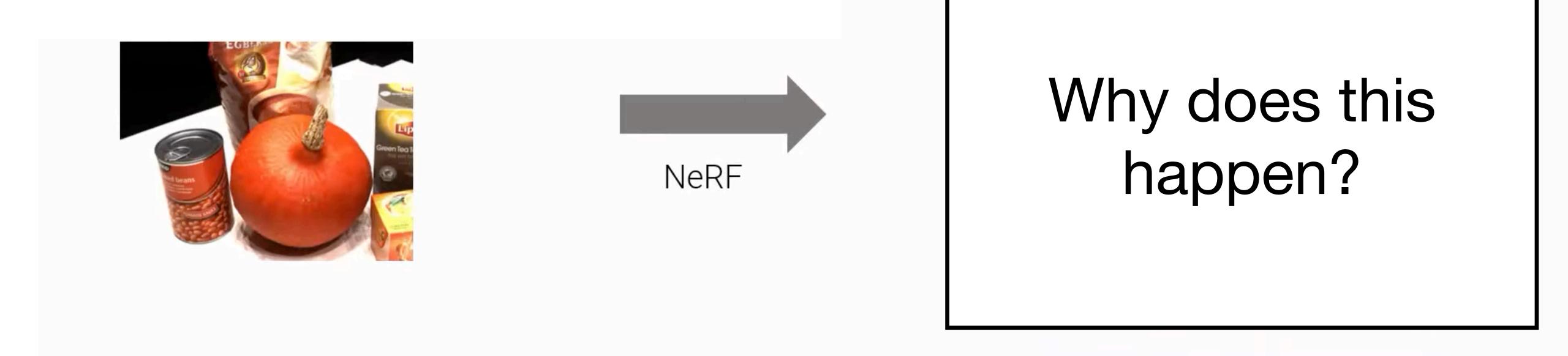
² Waymo

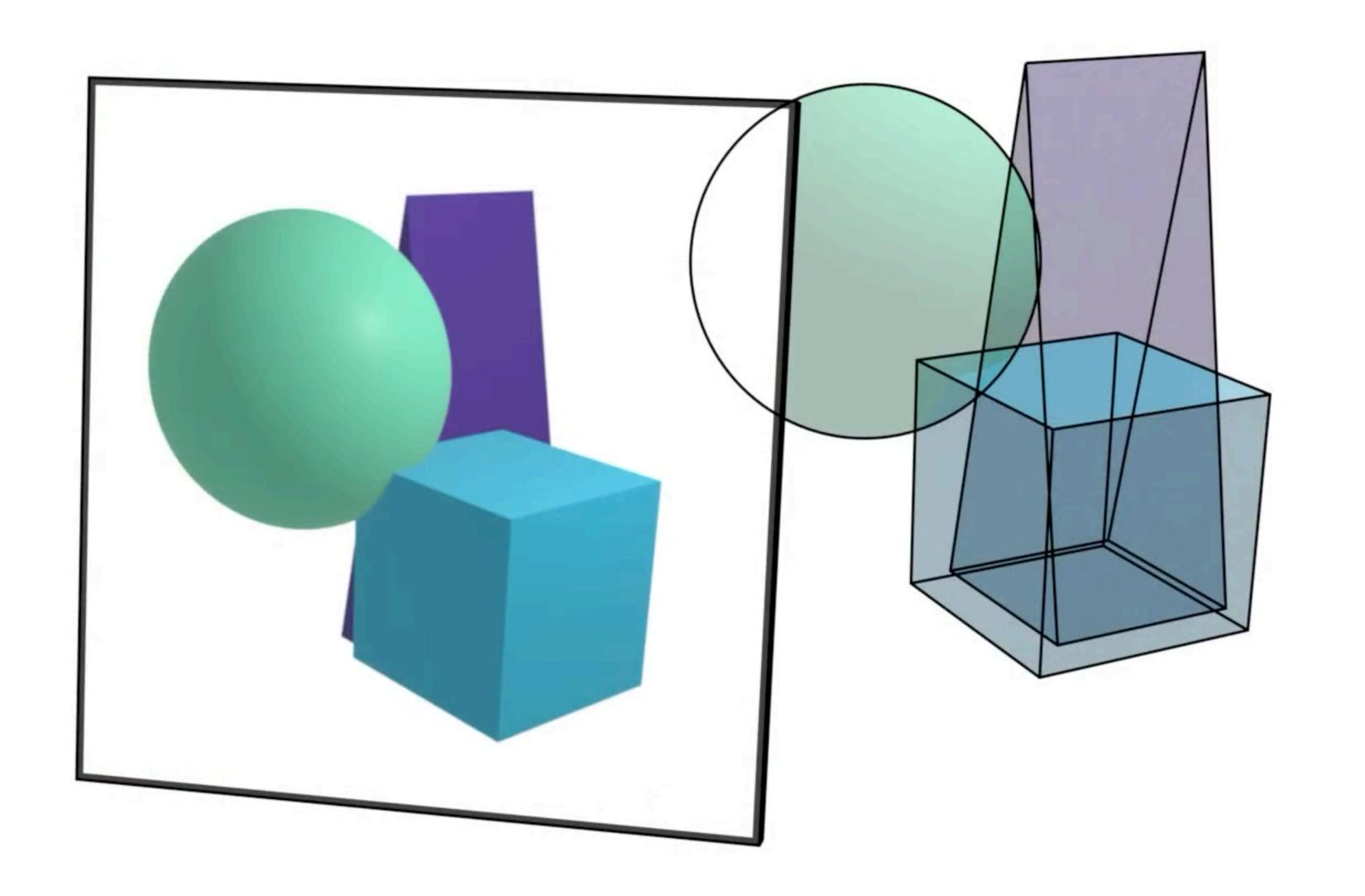
³Google Research



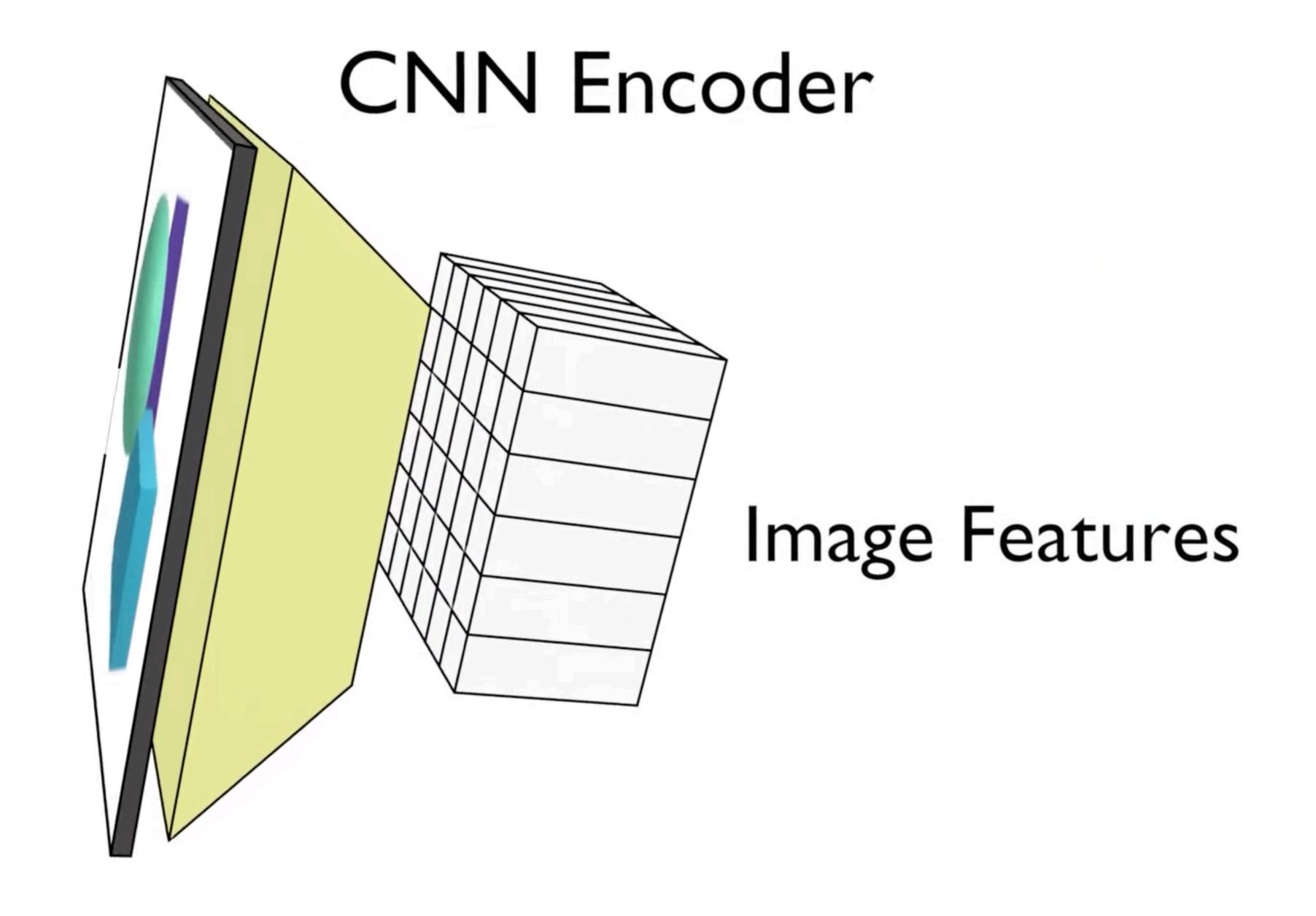


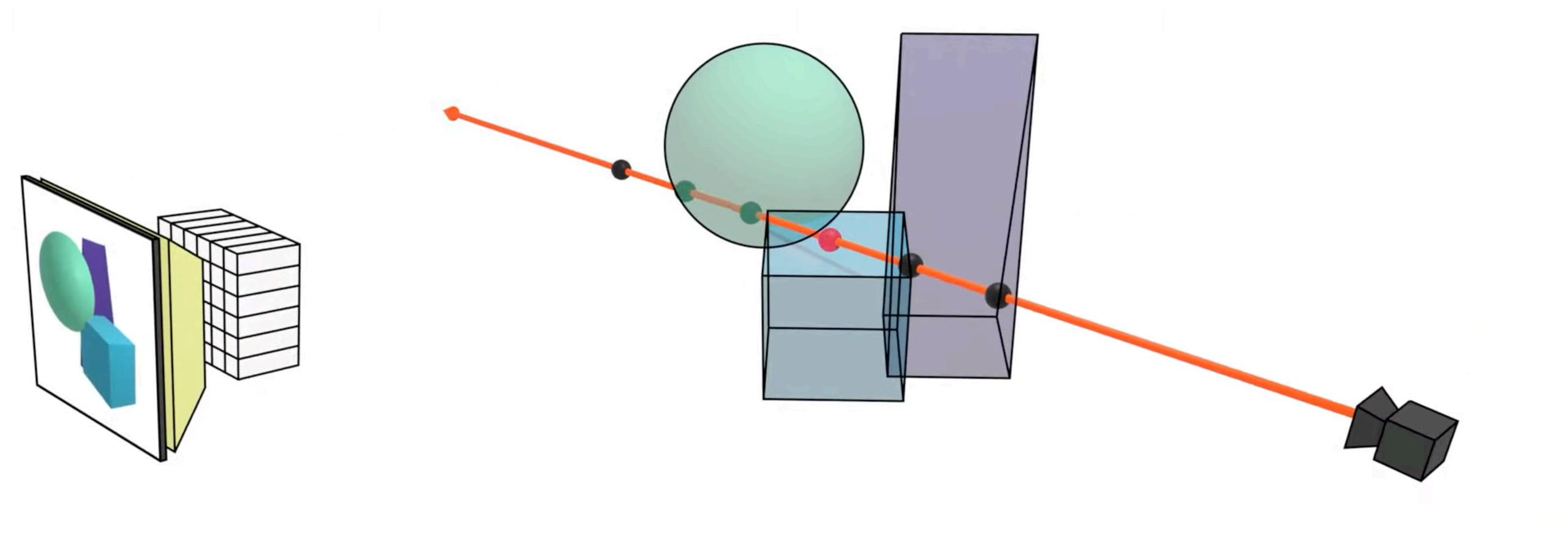
What happens if we use fewer views?

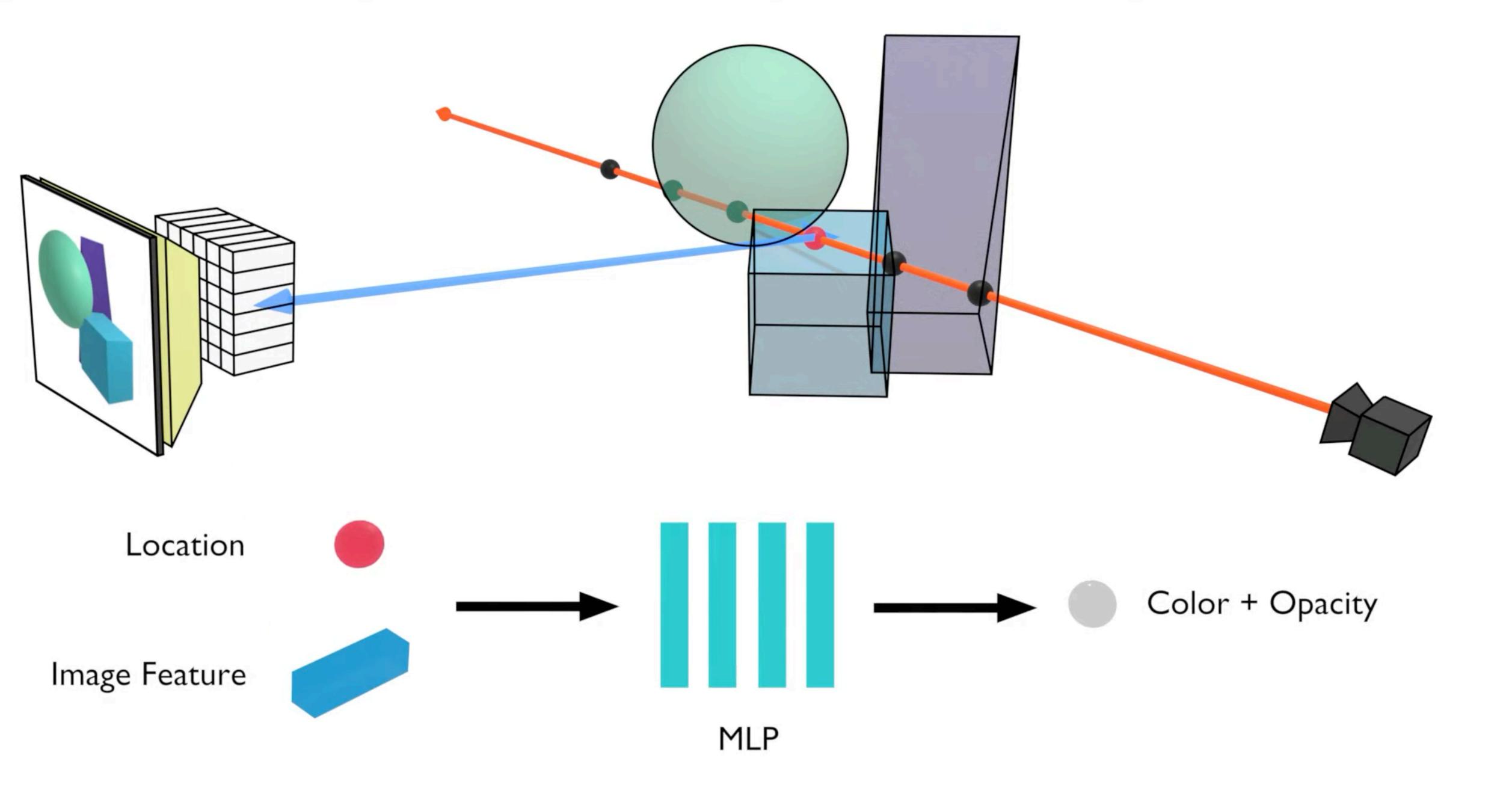




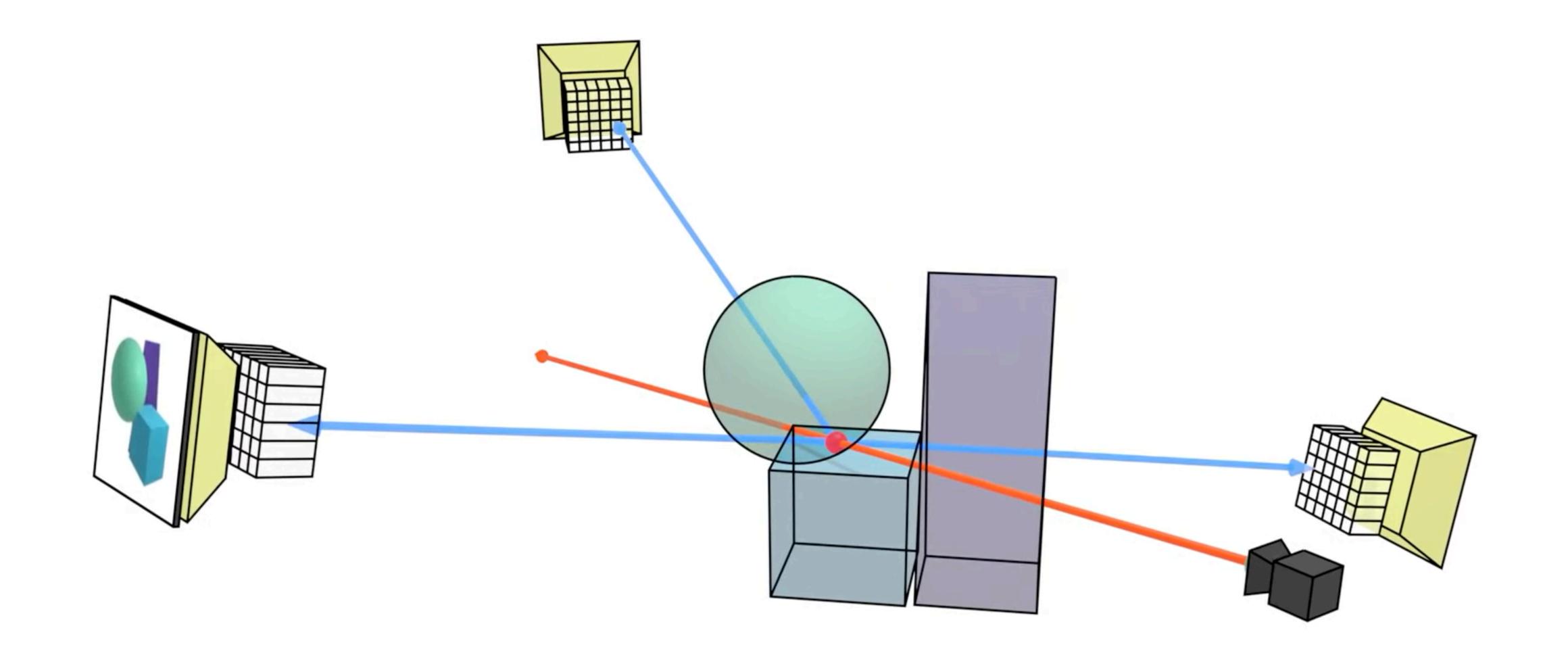
[Yu et al., "pixelNeRF: Neural Radiance Fields from One or Few Images", 2021]

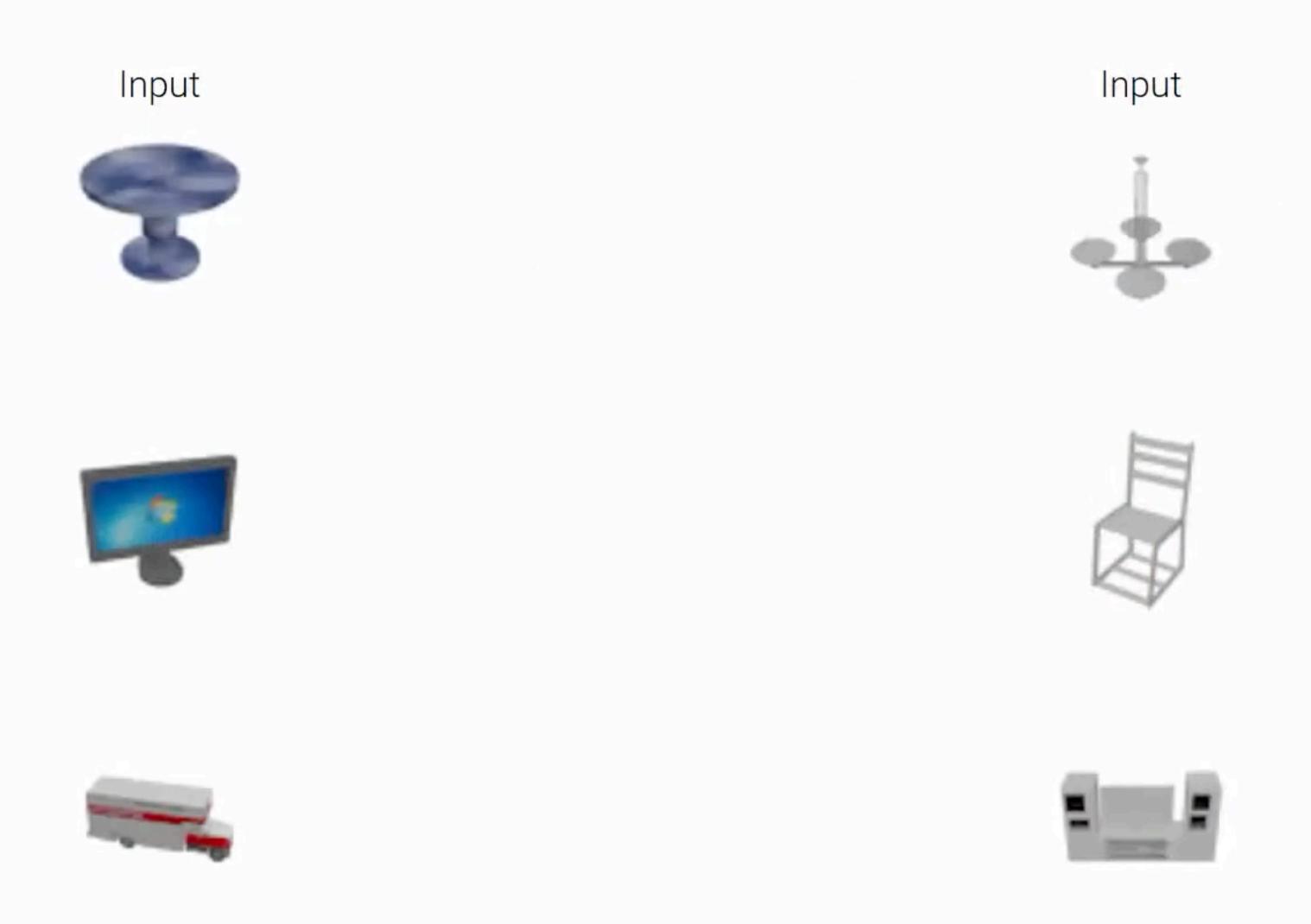






[Yu et al., "pixelNeRF: Neural Radiance Fields from One or Few Images", 2021]

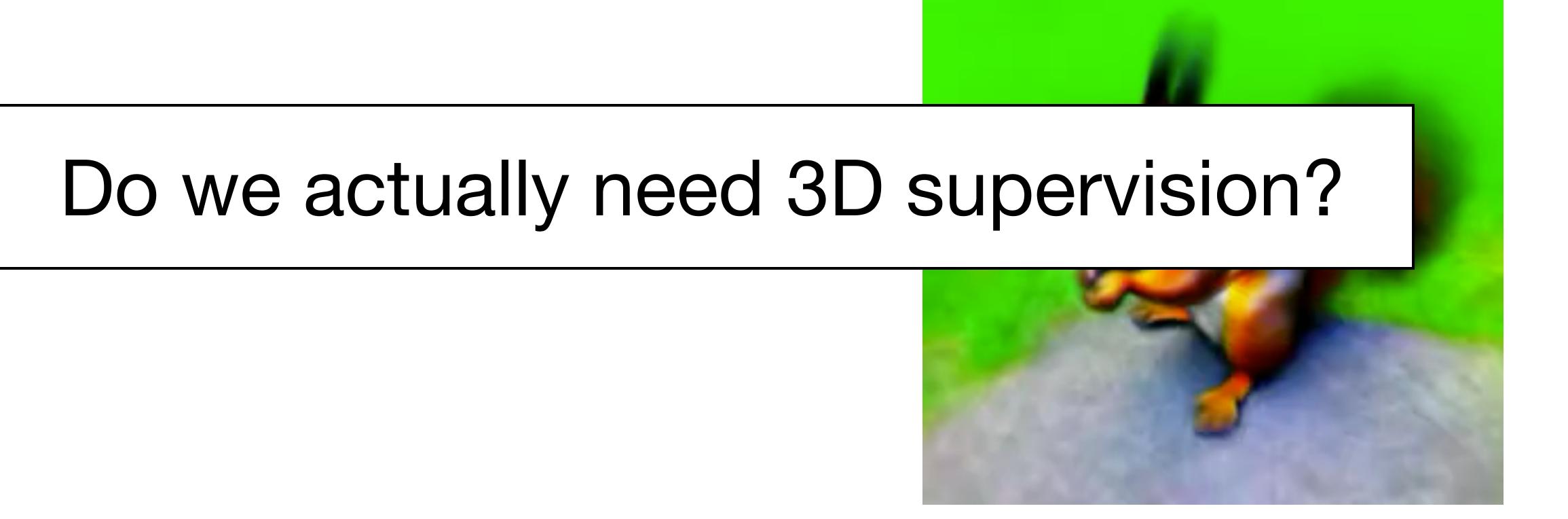




[Yu et al., "pixelNeRF: Neural Radiance Fields from One or Few Images", 2021]

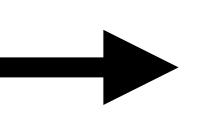
Can we generate 3D models?

Generating 3D models



We already have good text-to-image models

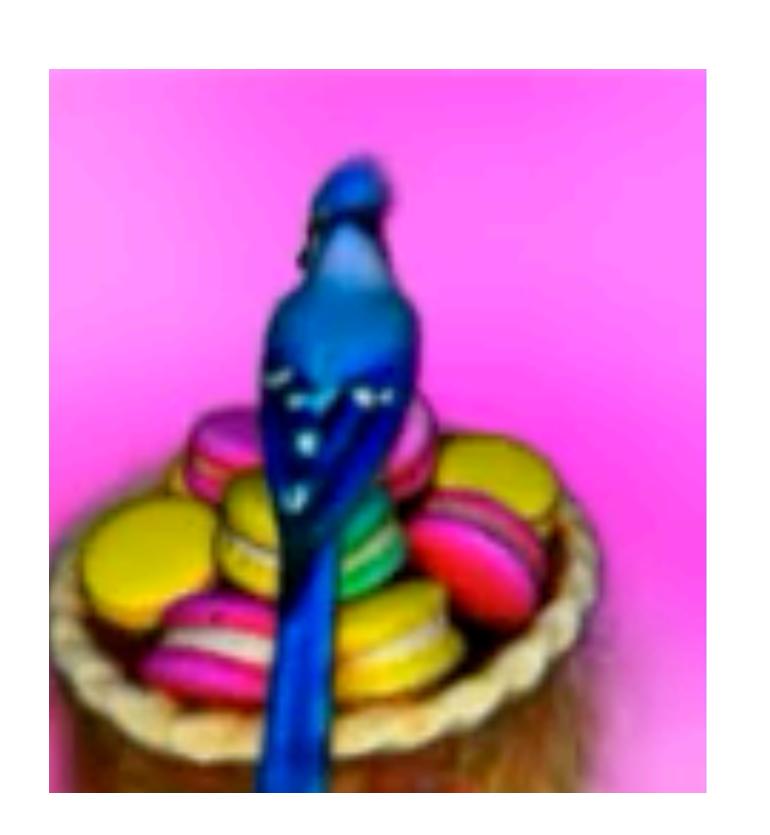
"A blue jay standing on a large basket of rainbow macarons."





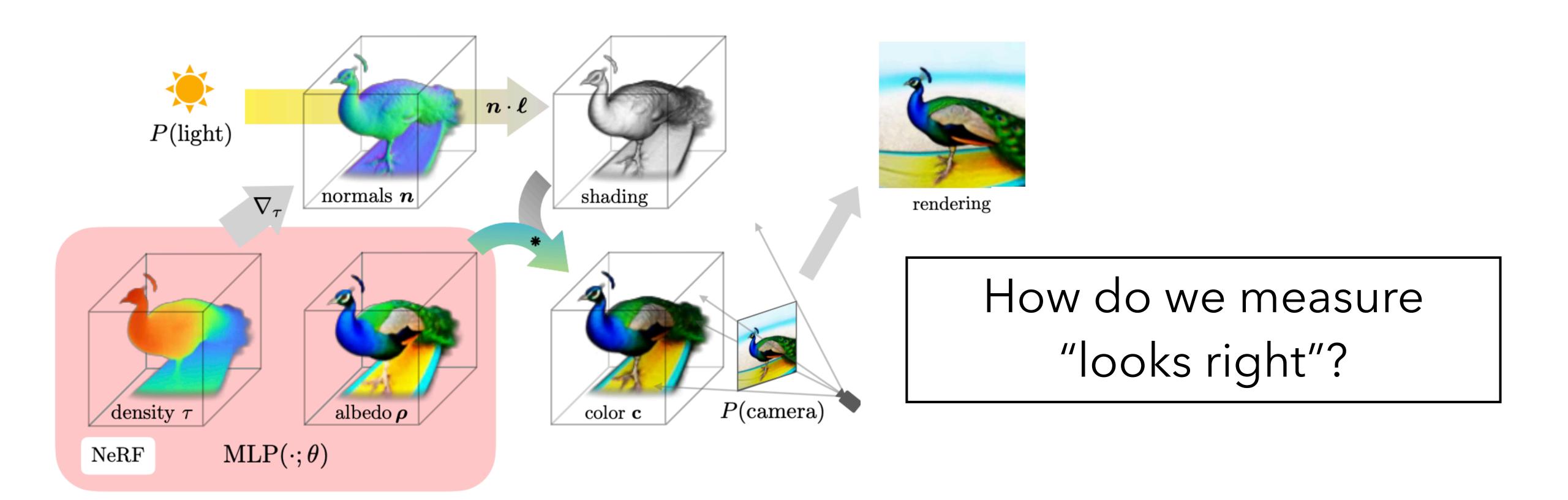
Could we fit a 3D model to this image?





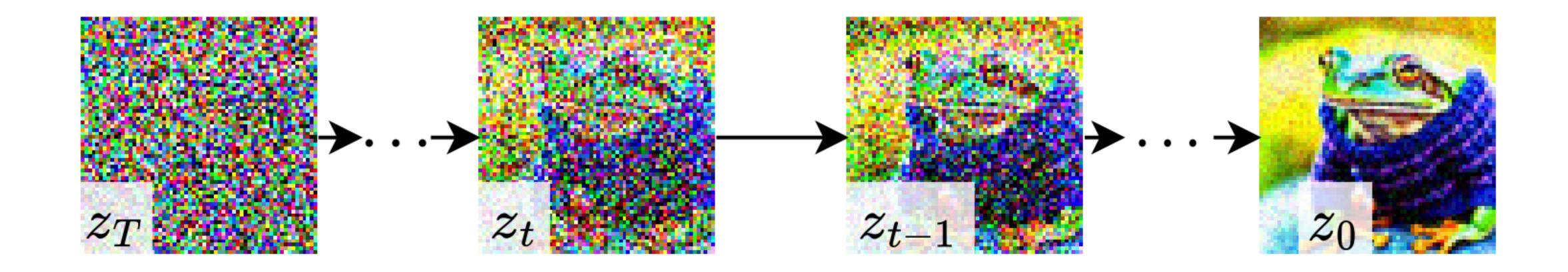
Single image could "supervise" single viewpoint.

DreamFusion: solve for a NeRF that "looks right" from every viewpoint



p(rendered image | "a DSLR photo of a peacock on a surfboard")

Can we use a pretrained diffusion model?

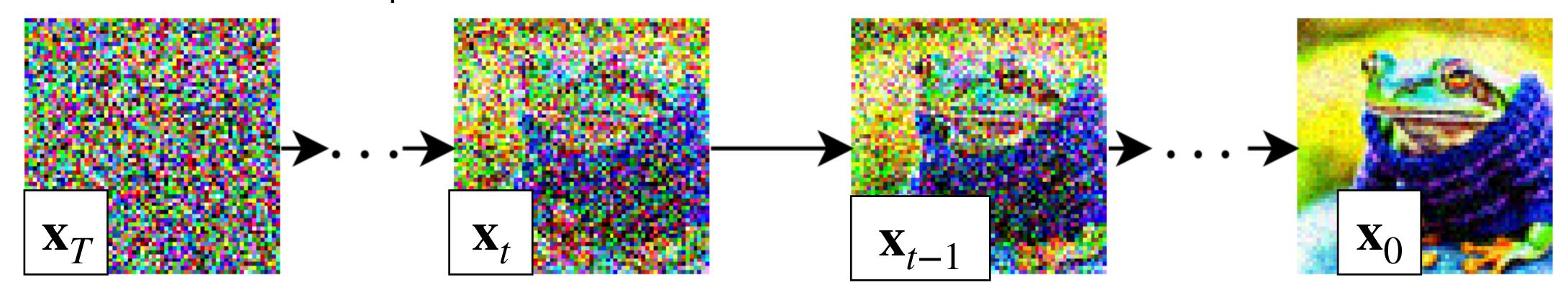


Diffusion: generate image by denoising

... but also it models *p*(image | text)!

Score distillation sampling

Reverse diffusion process:



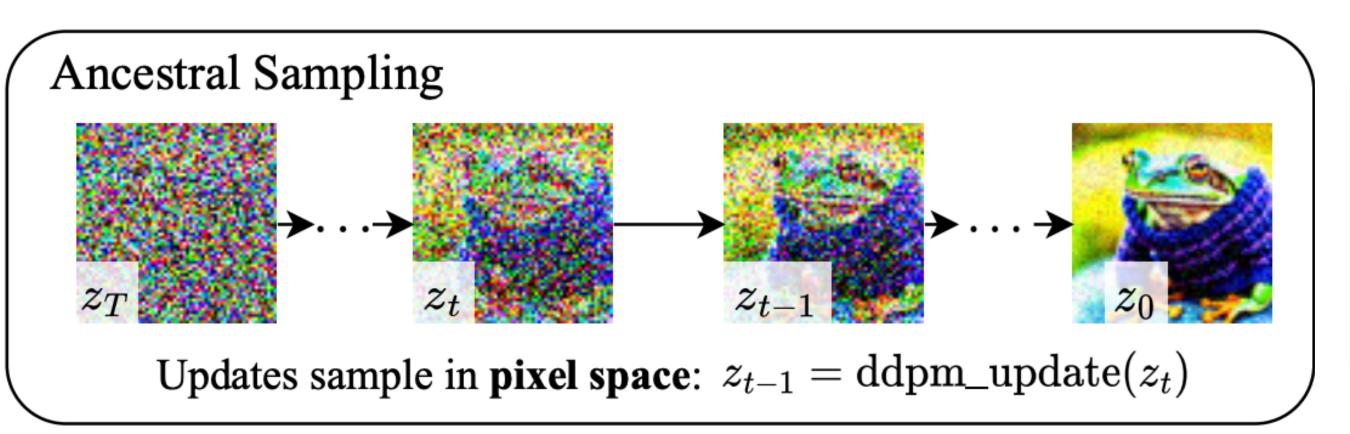
Diffusion objective function:

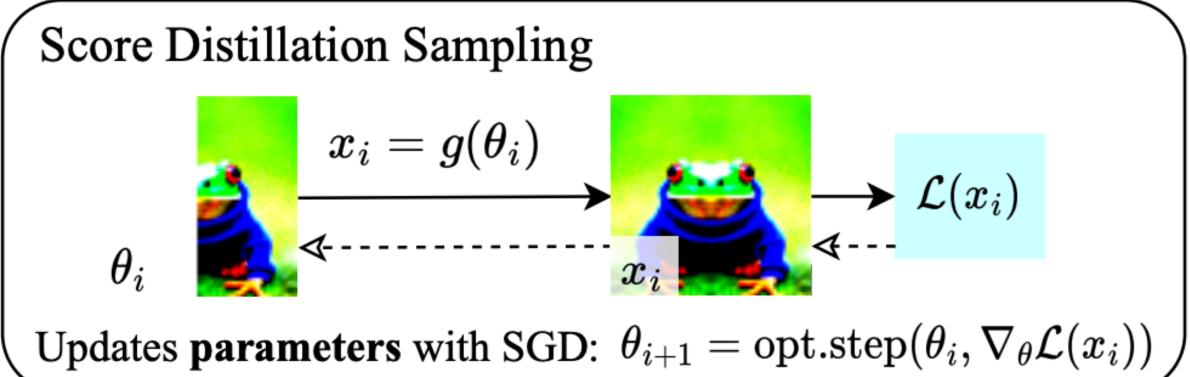
Diffusion objective function: Noise estimator Noise
$$\mathcal{L}_{\mathrm{Diff}}(\phi,\mathbf{x}) = \mathbb{E}_{t \sim \mathcal{U}(0,1)}, \quad [w(t)|\epsilon_{\phi}(\alpha_t\mathbf{x} + \sigma_t\epsilon;t) - \epsilon|_2^2]$$
 Noisy image

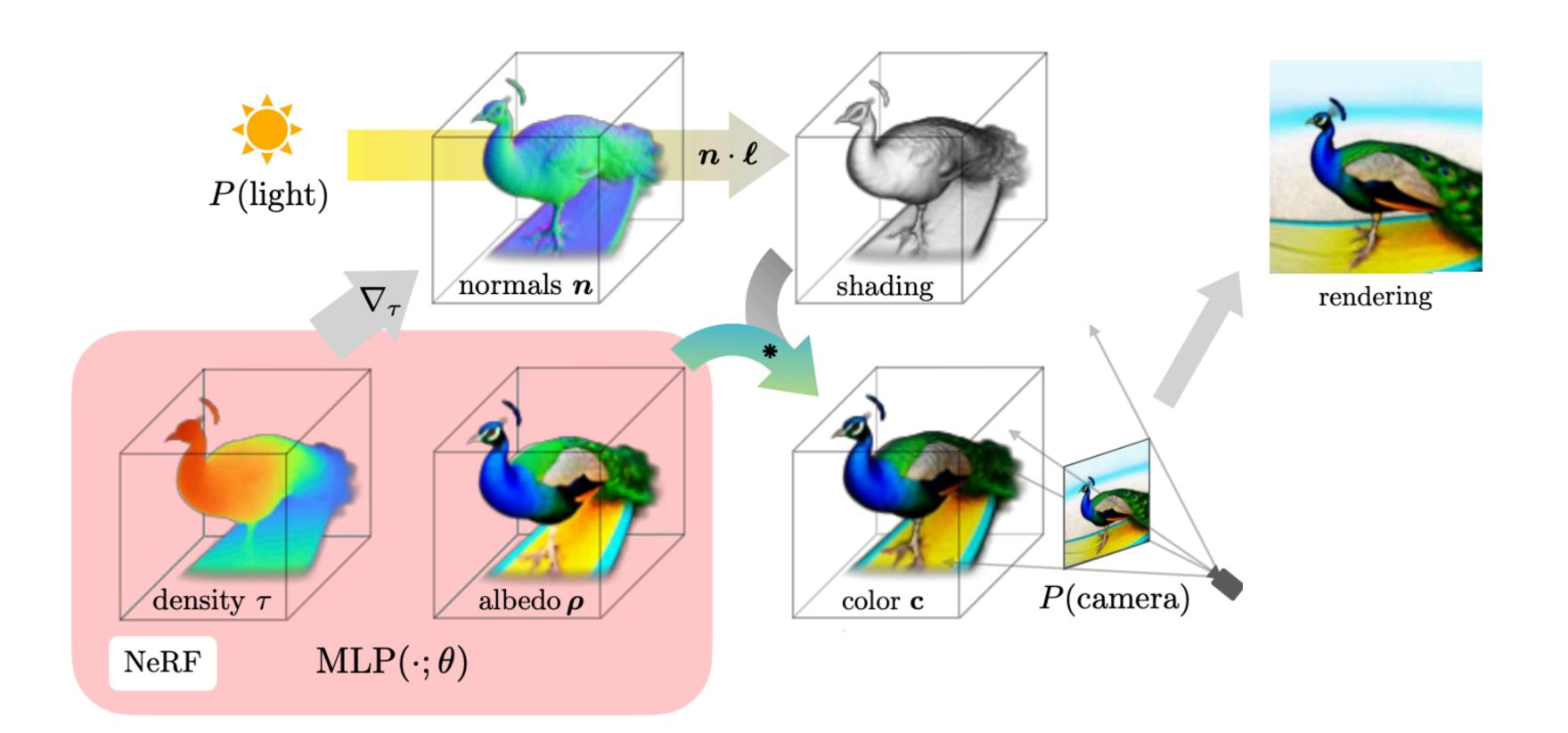
Noise estimator

Provides a lower bound on $log(p(\mathbf{x}))$ [Ho et al., 2020, Kingma et al., 2021]

Score distillation sampling







Solve for a NeRF such that, when it is rendered, it has high probability under a text-to-image model.

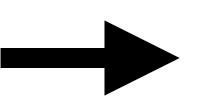
Source: [Poole et al., 2023]

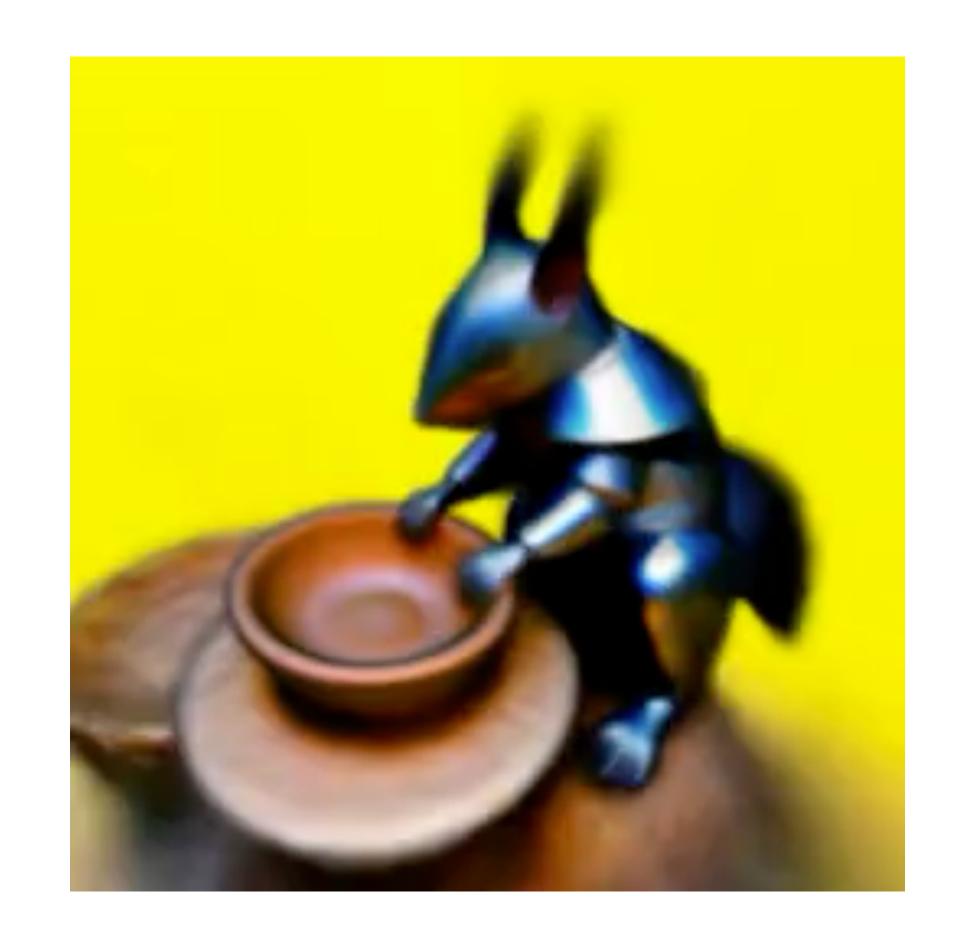
Generate 3D from text yourself!

```
a DSLR photo of a squirrel | an intricate wooden carving of a squirrel | a highly detailed metal sculpture of a squirrel
```

[...] | wearing a kimono | wearing a medieval suit of armor | wearing a purple hoodie | wearing an elegant ballgown

[...] | reading a book | riding a motorcycle | playing the saxophone | chopping vegetables | sitting at a pottery wheel shaping a clay bowl | riding a skateboard | wielding a katana | eating a hamburger | dancing





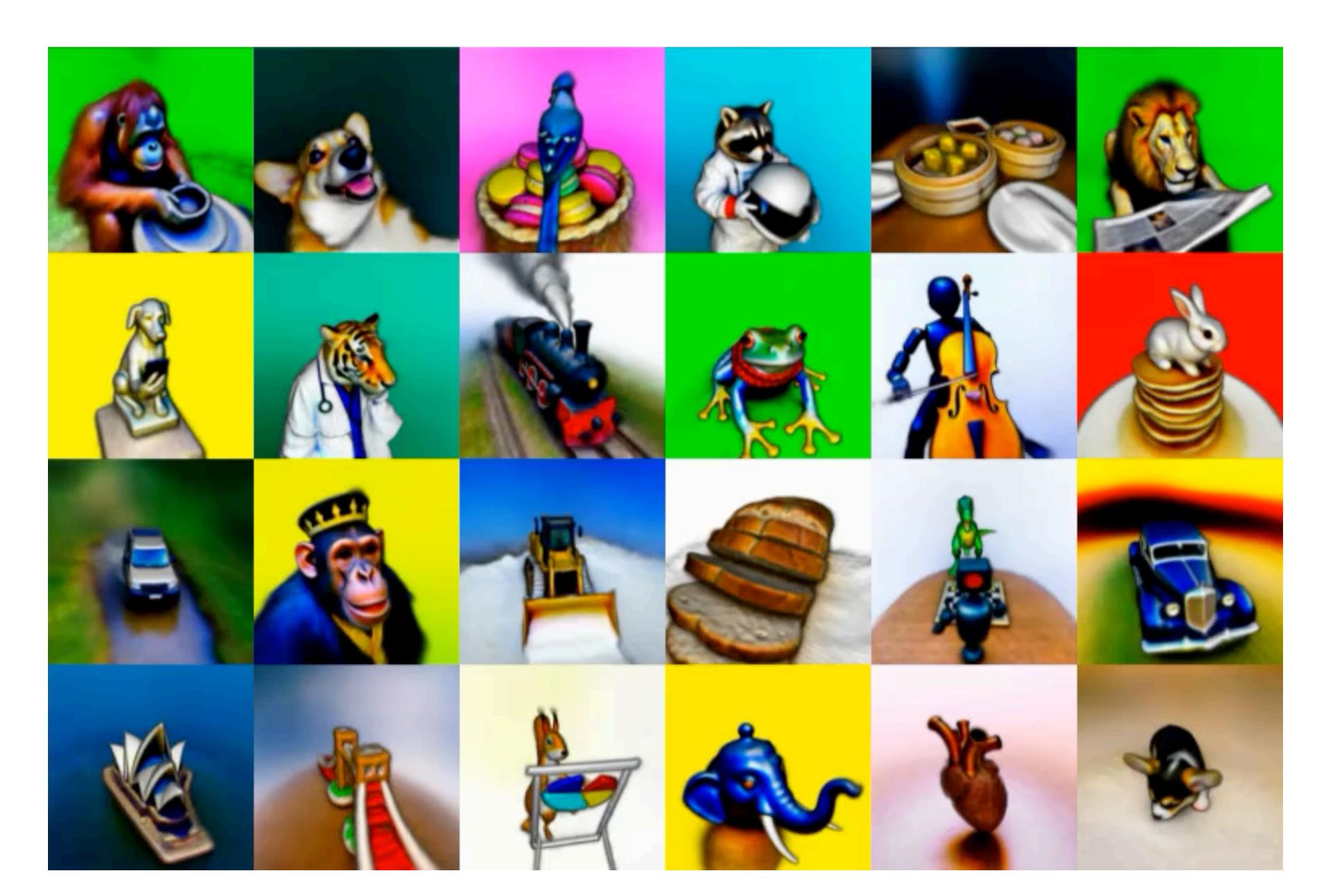
Generate 3D from text yourself!

```
a DSLR photo of a squirrel | an intricate wooden carving of a squirrel | a highly detailed metal sculpture of a squirrel
```

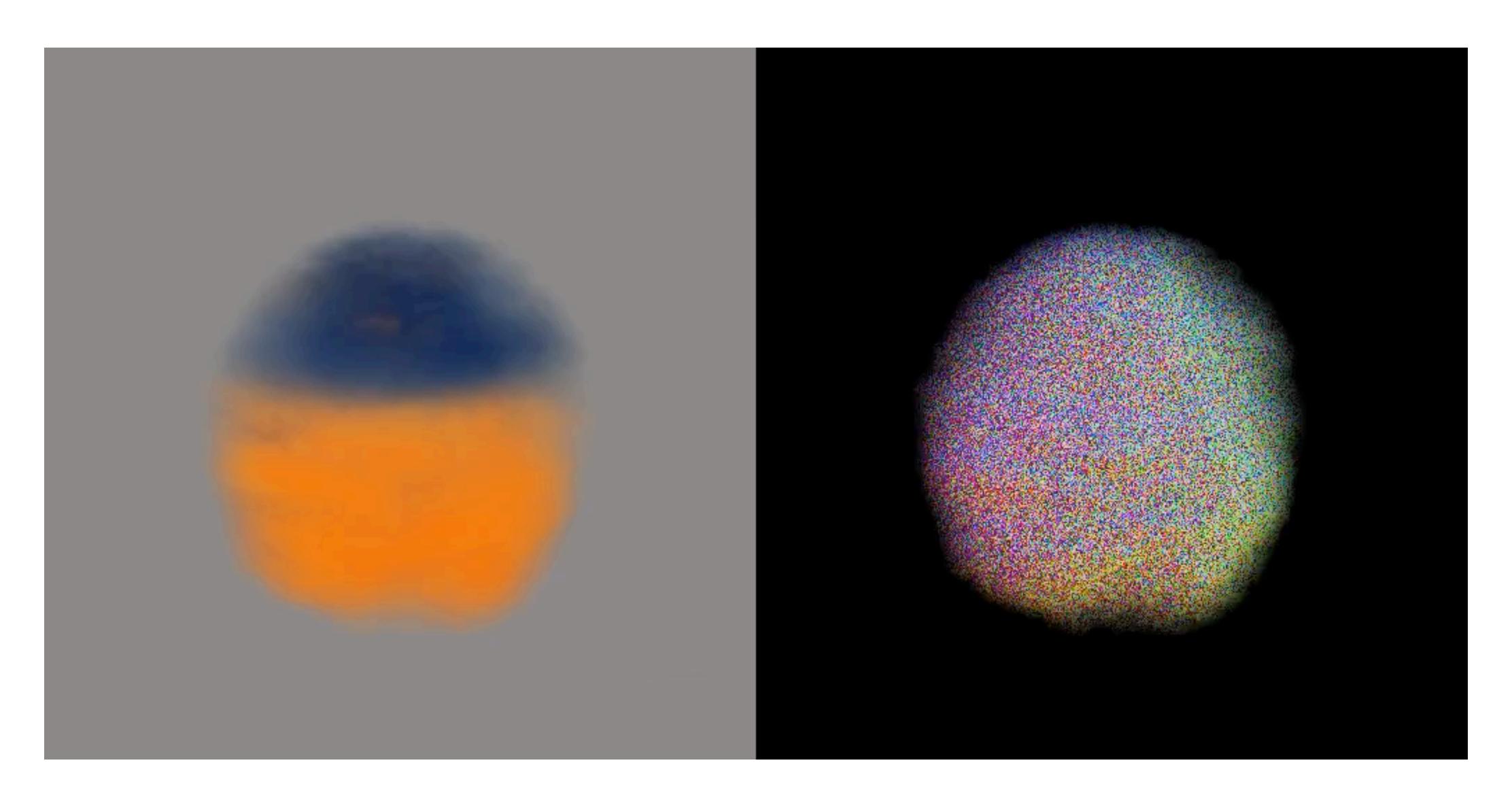
[...] | wearing a kimono | wearing a medieval suit of armor | wearing a purple hoodie | wearing an elegant ballgown

[...] | reading a book | riding a motorcycle | playing the saxophone | chopping vegetables | sitting at a pottery wheel shaping a clay bowl | riding a skateboard | wielding a katana | eating a hamburger | dancing

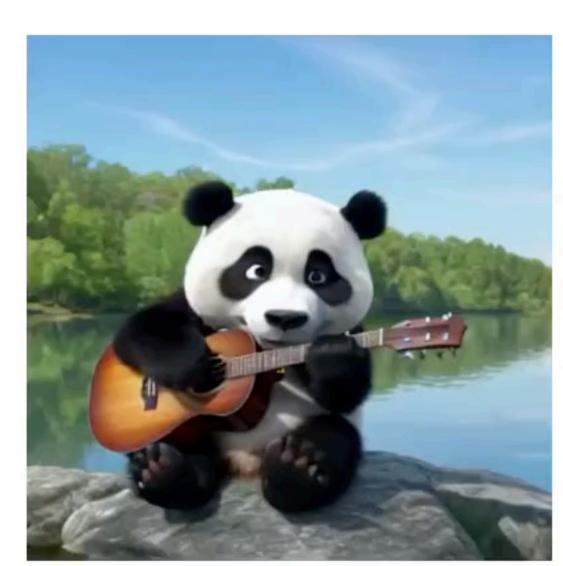




More recent text-to-3D

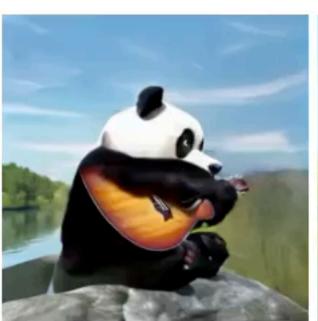


Text-to-4D



Input video

Sample from multi-view video diffusion model

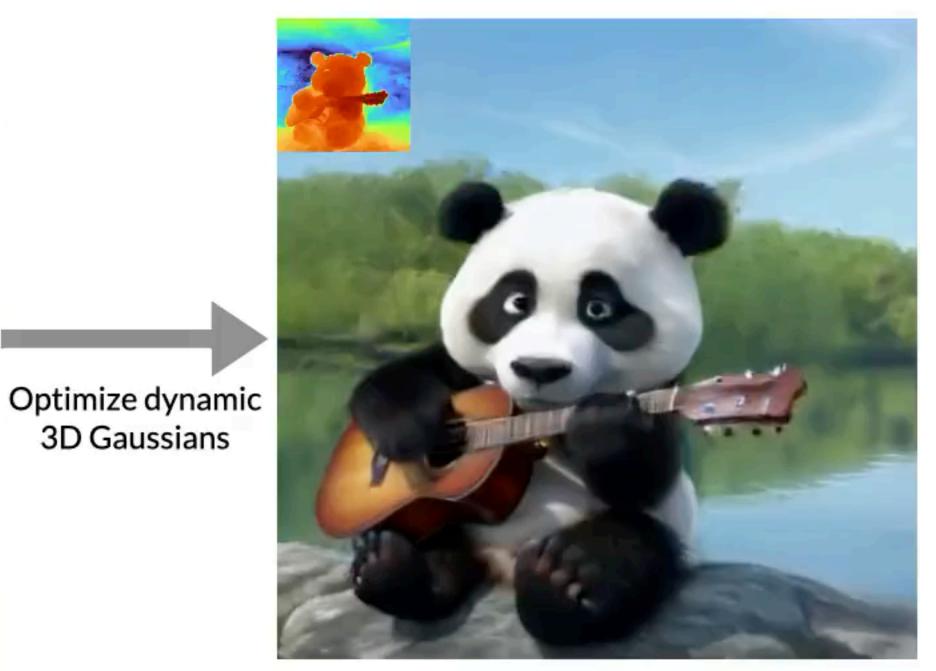




3D Gaussians



Generated multi-view videos



Dynamic 3D scene



[Wu et al., CAT4D, 2024]

Next class: light