

Neural Radiance Fields pt 2



Video from the original ECCV'20 paper

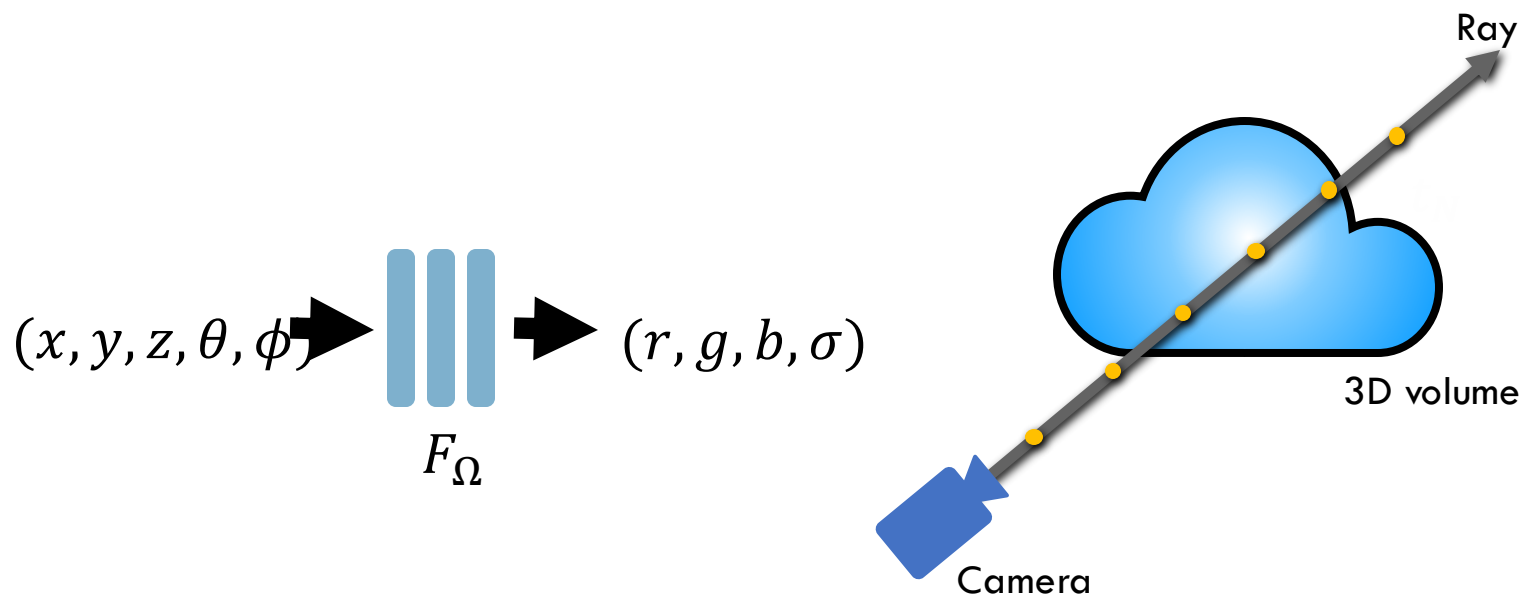
CS180/280A: Intro to Computer Vision and Computational Photography

Angjoo Kanazawa and Alexei Efros

UC Berkeley Fall 2025

Lots of content from Noah Snavely and [ECCV 2022 Tutorial on Neural Volumetric Rendering for Computer Vision](#)

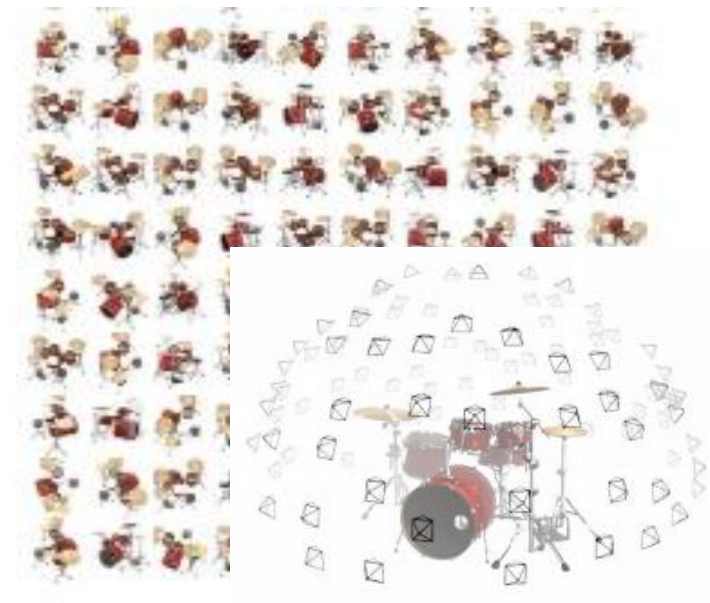
Recap: 3 Key Components



Neural Volumetric 3D
Scene Representation

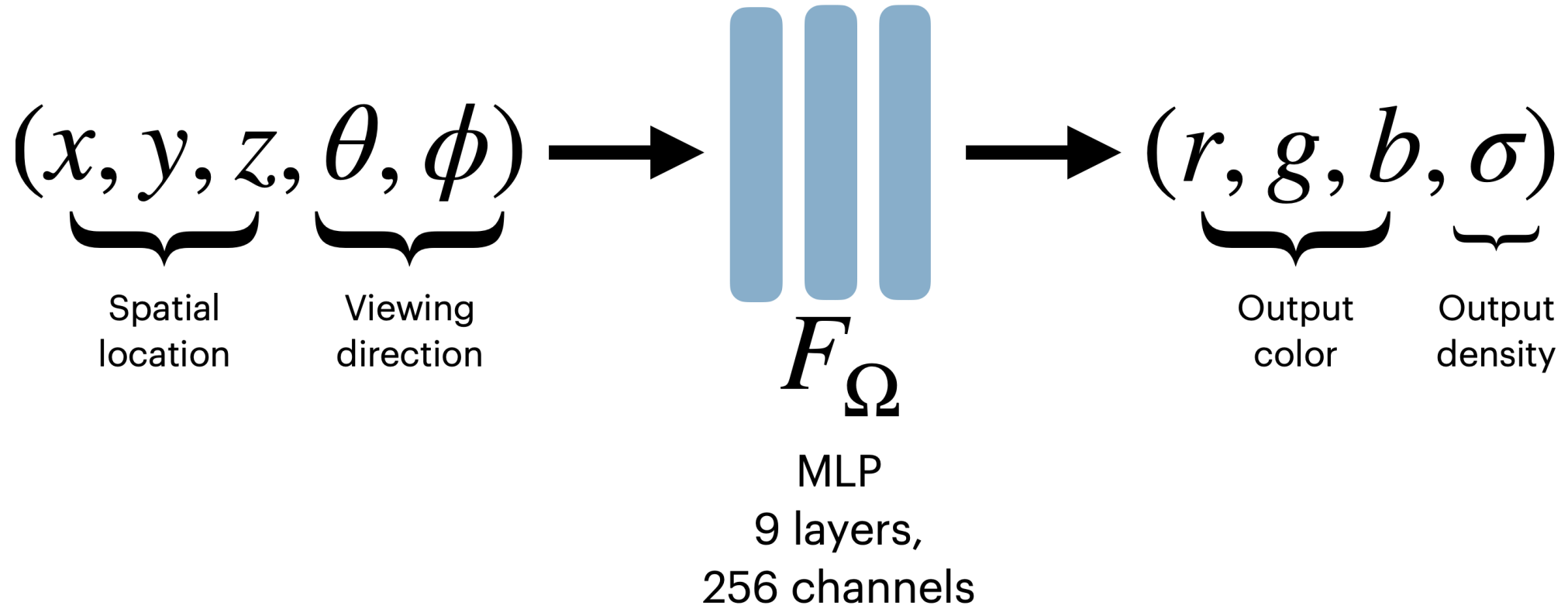
Differentiable Volumetric
Rendering Function

Objective: Synthesize
all training views



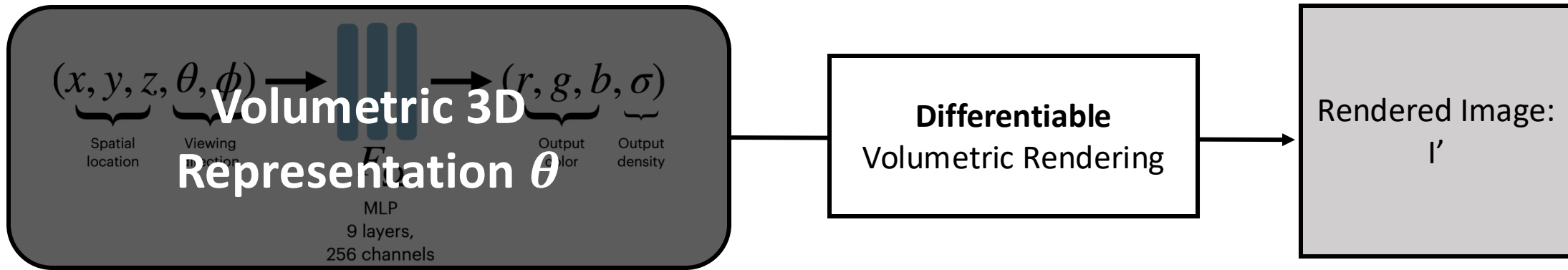
Optimization via
Analysis-by-Synthesis

“Neural Radiance Fields”



“Neural Radiance Fields”

How an image is made (“Inference”)



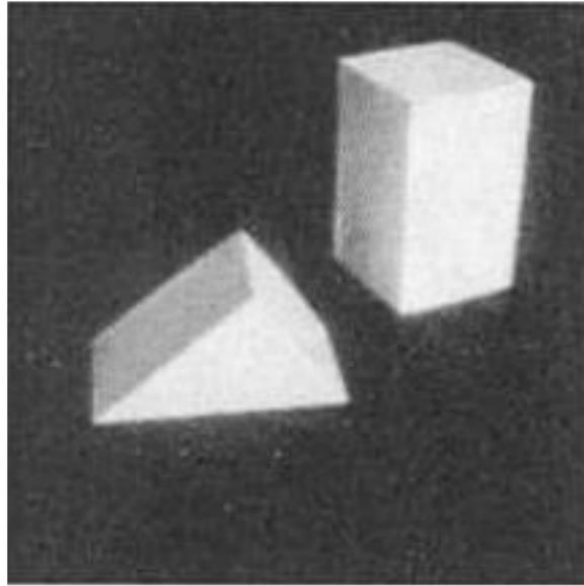
“Training” Objective (aka Analysis-by-Synthesis):

$$\min_{\theta} \| \text{Rendered Image: } I' - \text{Observed Image: } I \|_2$$

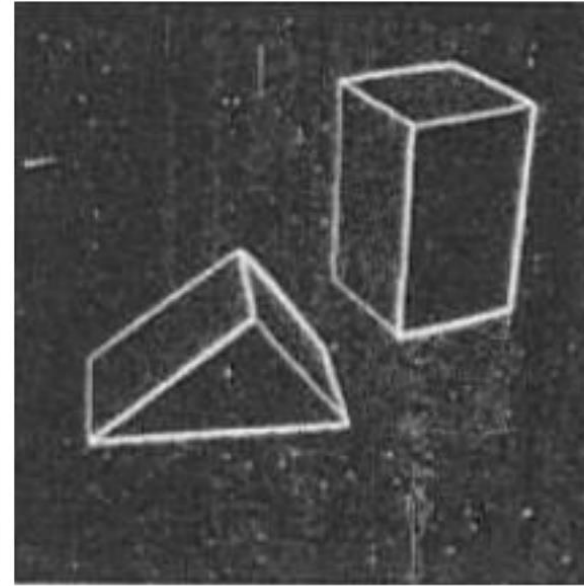
Analysis-by-Synthesis



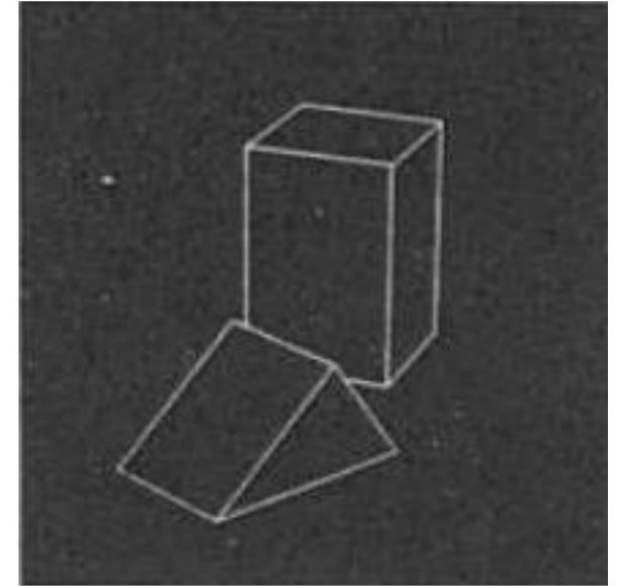
Larry Roberts
“Father of Computer Vision”



Input image



2x2 gradient operator

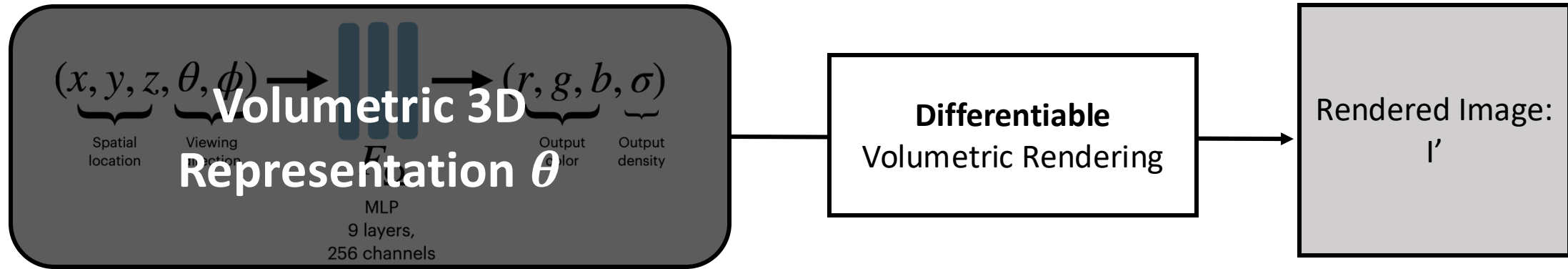


computed 3D model
rendered from new viewpoint

- History goes way back to the **first** Computer Vision paper!
Roberts: *Machine Perception of Three-Dimensional Solids*, MIT, 1963

“Neural Radiance Fields”

Forward Function: How an image is made (Inference)



“Training” Objective (aka Analysis-by-Synthesis):

$$\min_{\theta} \| \text{Rendered Image: } I' - \text{Observed Image: } I \|_2$$

Differentiable Rendering

- How to change θ (network parameter) so that we get the final image?
- Gradient Descent “Hiking”

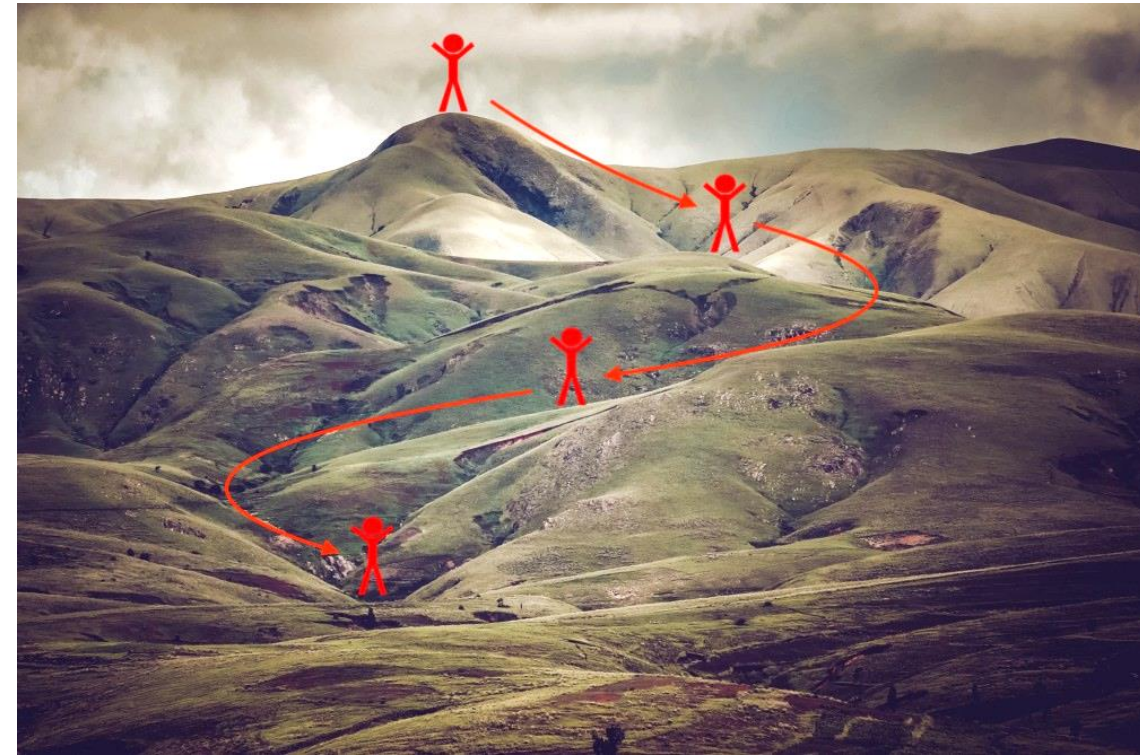
Same idea here, “hiking” now means you’re going to change the network parameter little by little.

The “Mountain” or the “Loss” comes from the reconstruction loss.

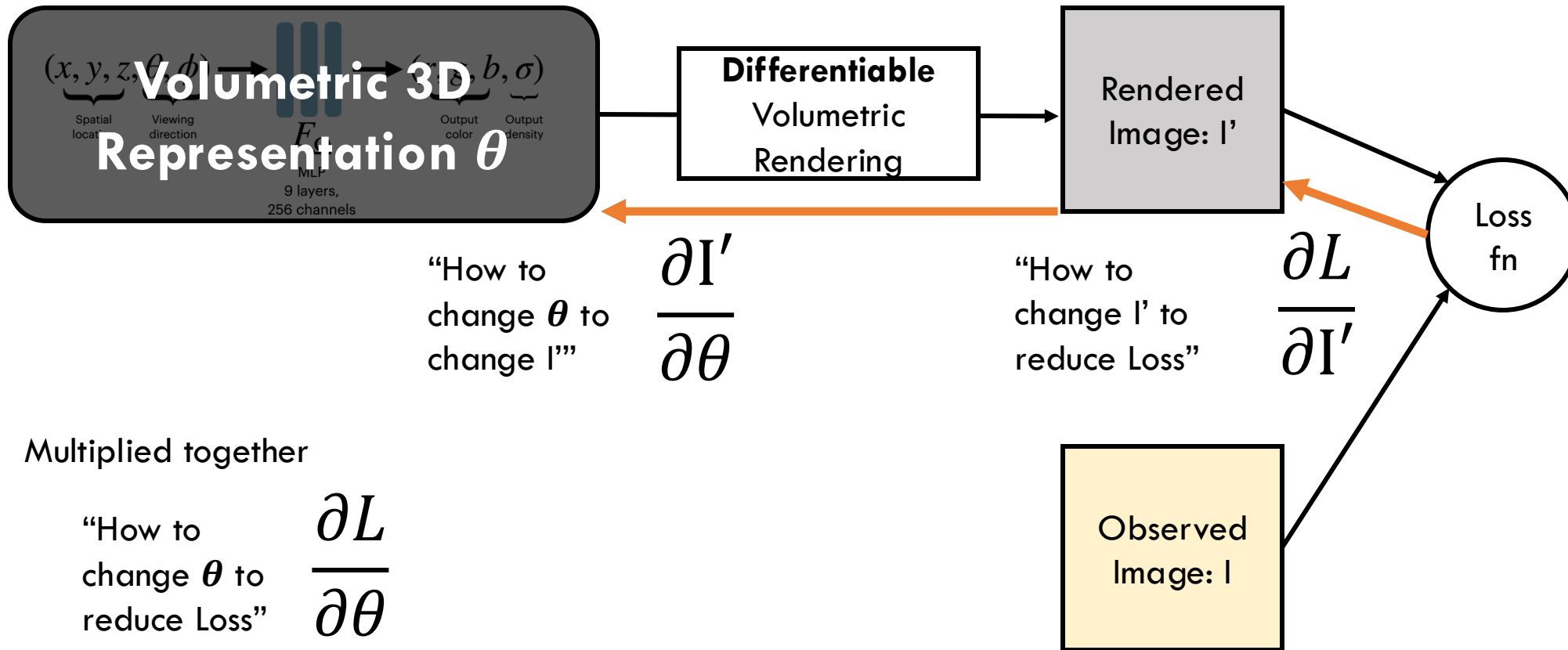
$$\frac{\partial L}{\partial \theta} = \frac{\partial L}{\partial I'} \frac{\partial I'}{\partial \theta}$$

$$L = \|I' - I\|$$
$$I' = f(x; \theta)$$

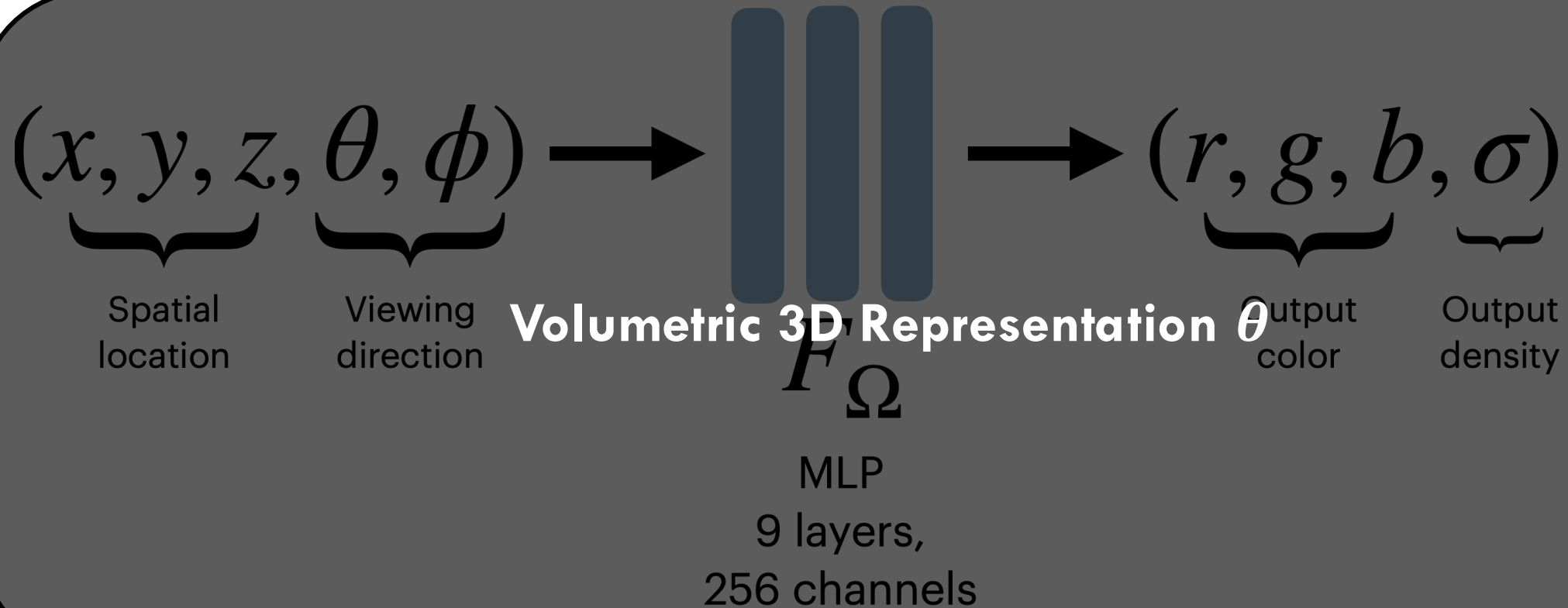
Chain rule, aka Back propagation



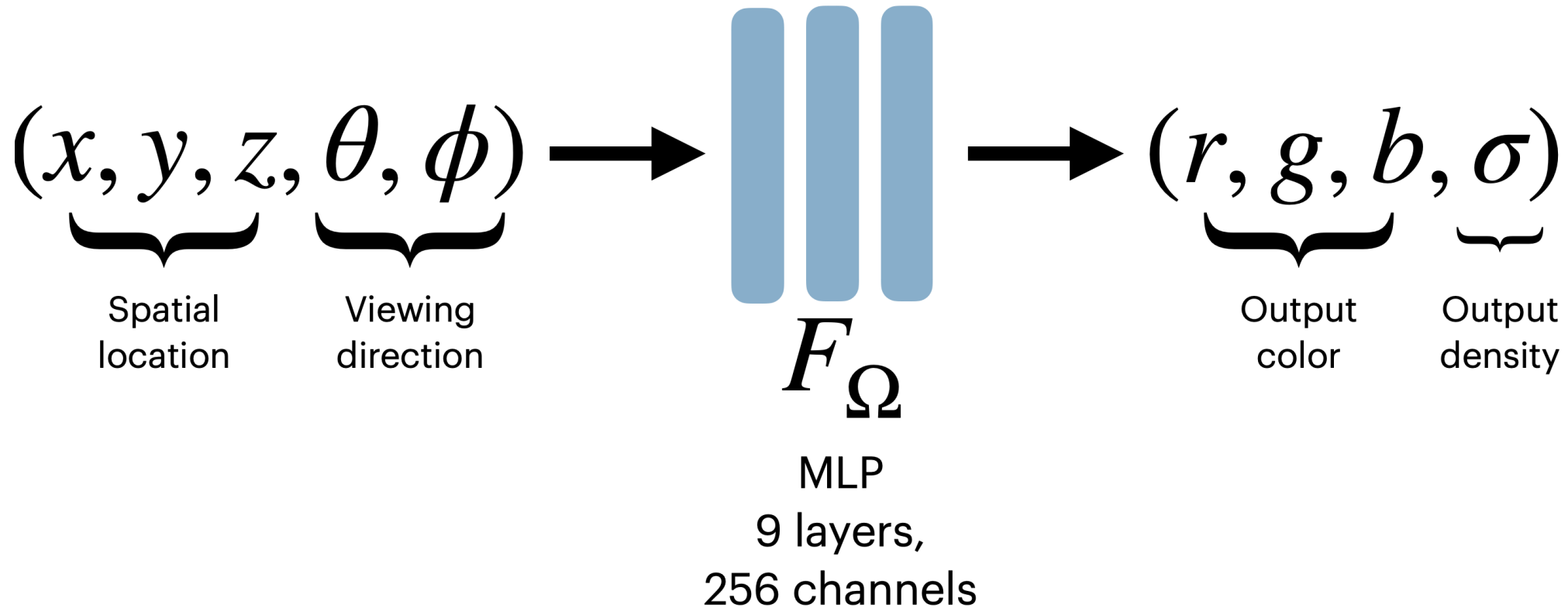
“Neural Radiance Fields”



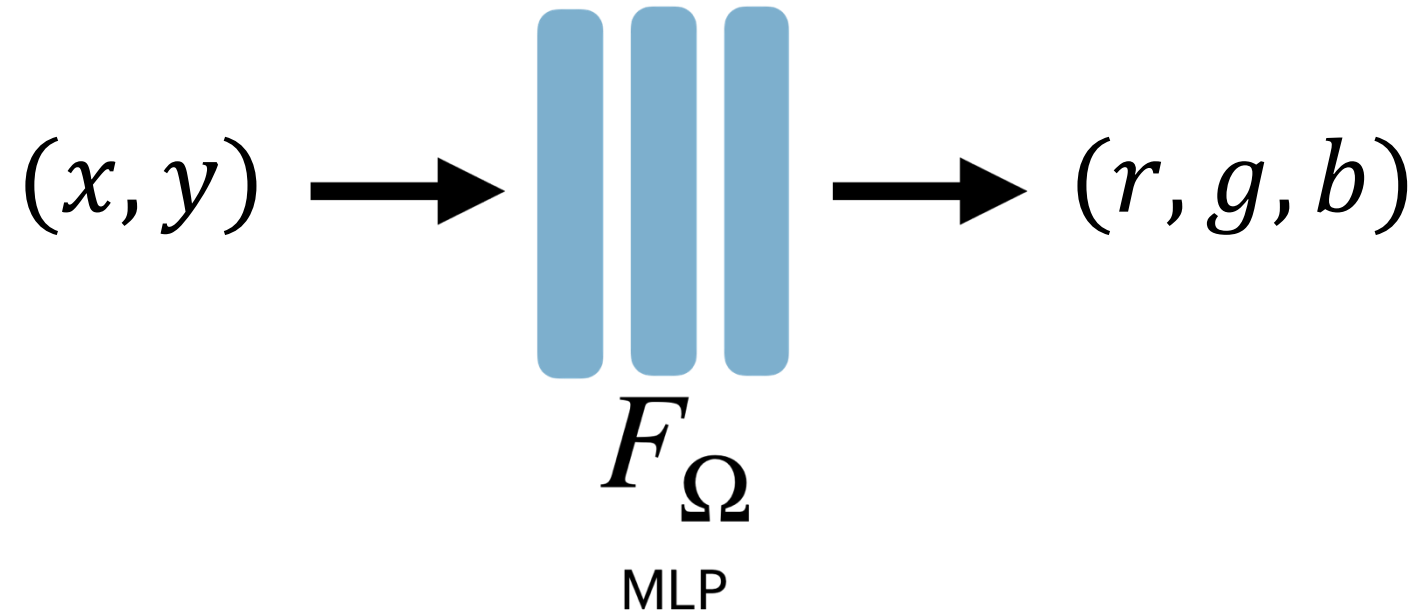
“Neural Radiance Fields”



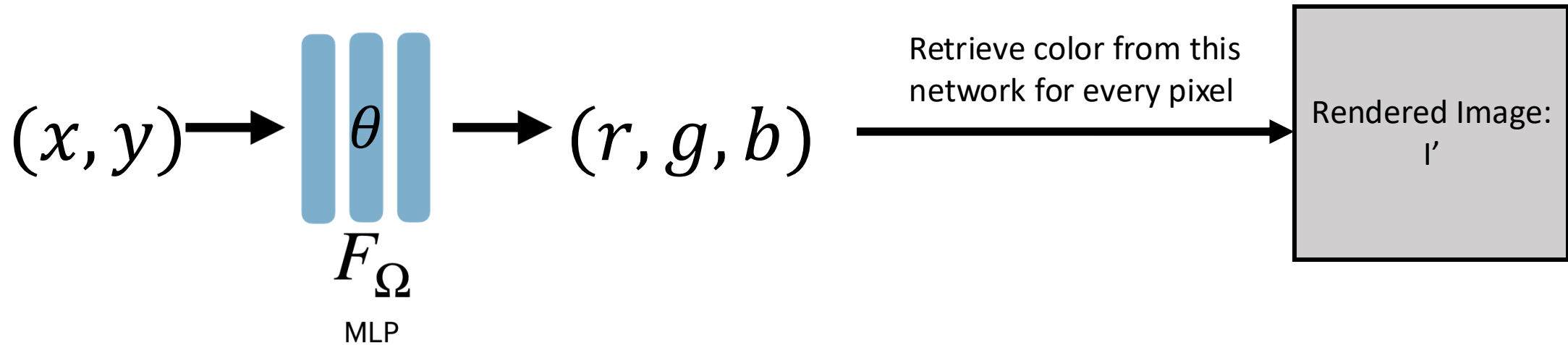
“Neural Radiance Fields”



Let's simplify, do this in 2D:



Let's simplify, do this in 2D:

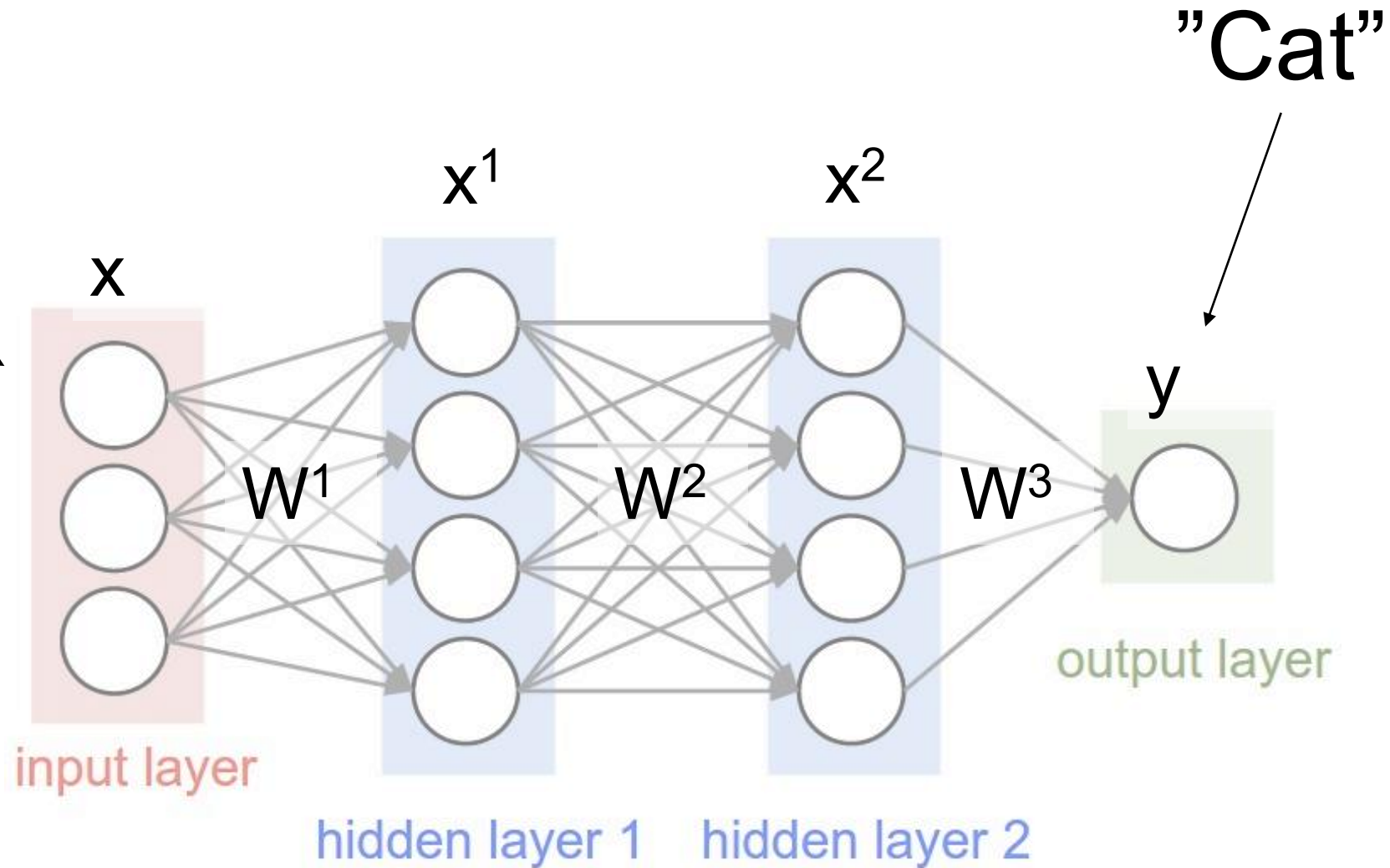


Optimize with “Training” Objective (aka Analysis-by-Synthesis):

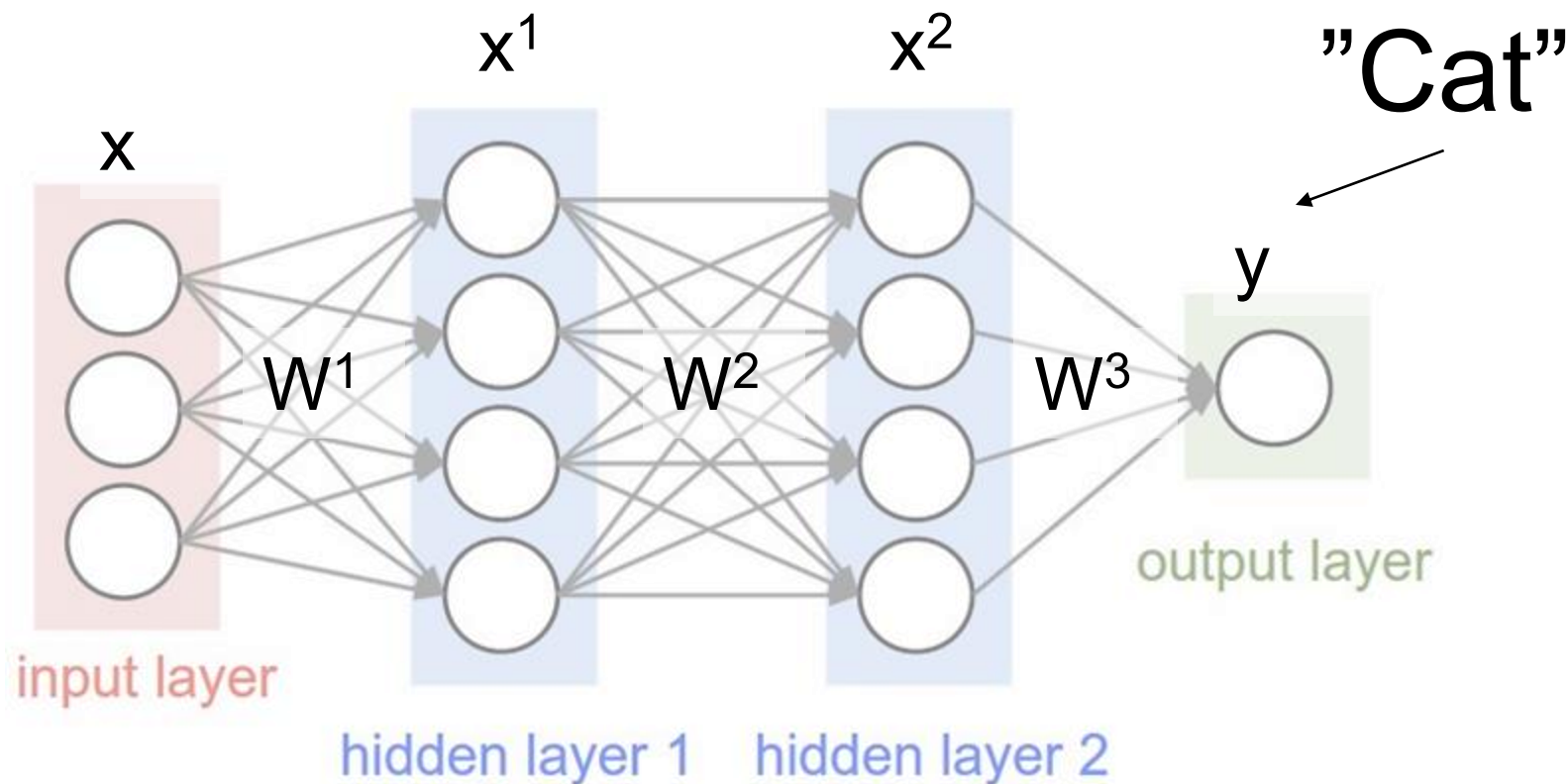
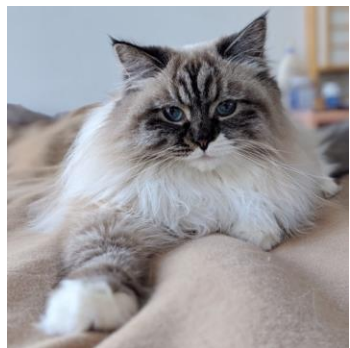
$$\frac{\partial L}{\partial \theta} = \frac{\partial (rgb - rgb')}{\partial \theta} \quad \min_{\theta} \left\| \begin{array}{|c|} \hline \text{Rendered Image: } I' \\ \hline \end{array} - \begin{array}{|c|} \hline \text{Observed Image: } I \\ \hline \end{array} \right\|_2$$

Straight forward to implement with Pytorch

ML Recap: Multi-layer perceptrons / Fully-Connected Layer



Multi-layer perceptrons / Fully-Connected Layer



In each layer:

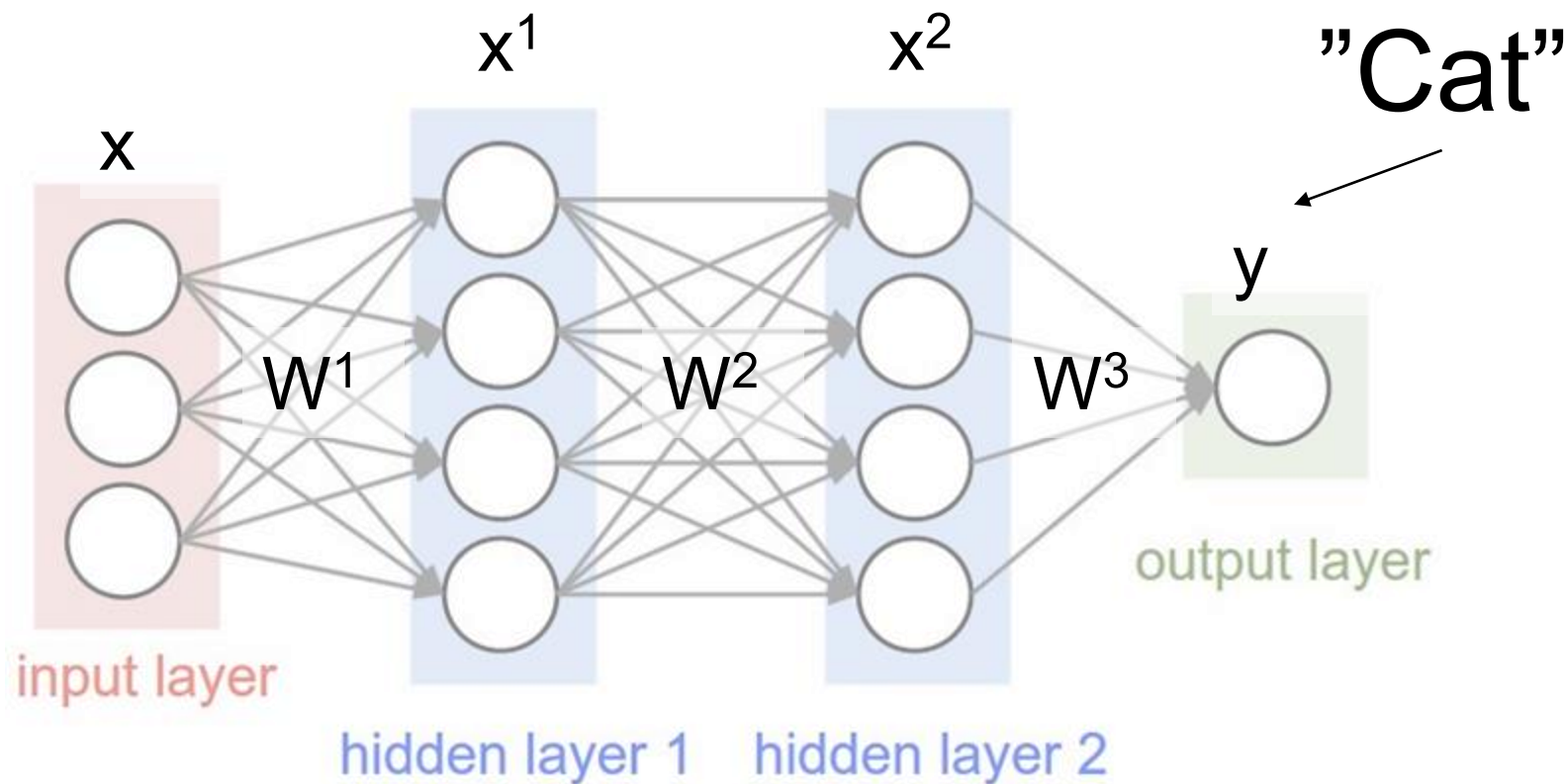
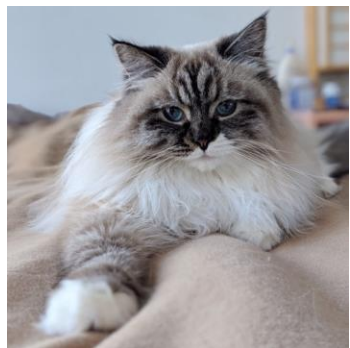
1. Linear Transform $z = W^l x^{l-1} + b$
2. Apply Non-Linearity $x^l = f(z)$

Usually

$$f = \text{RELU}(z) \\ = \max(0, z)$$

what
happens if f
is identity?

Multi-layer perceptrons / Fully-Connected Layer



In each layer:

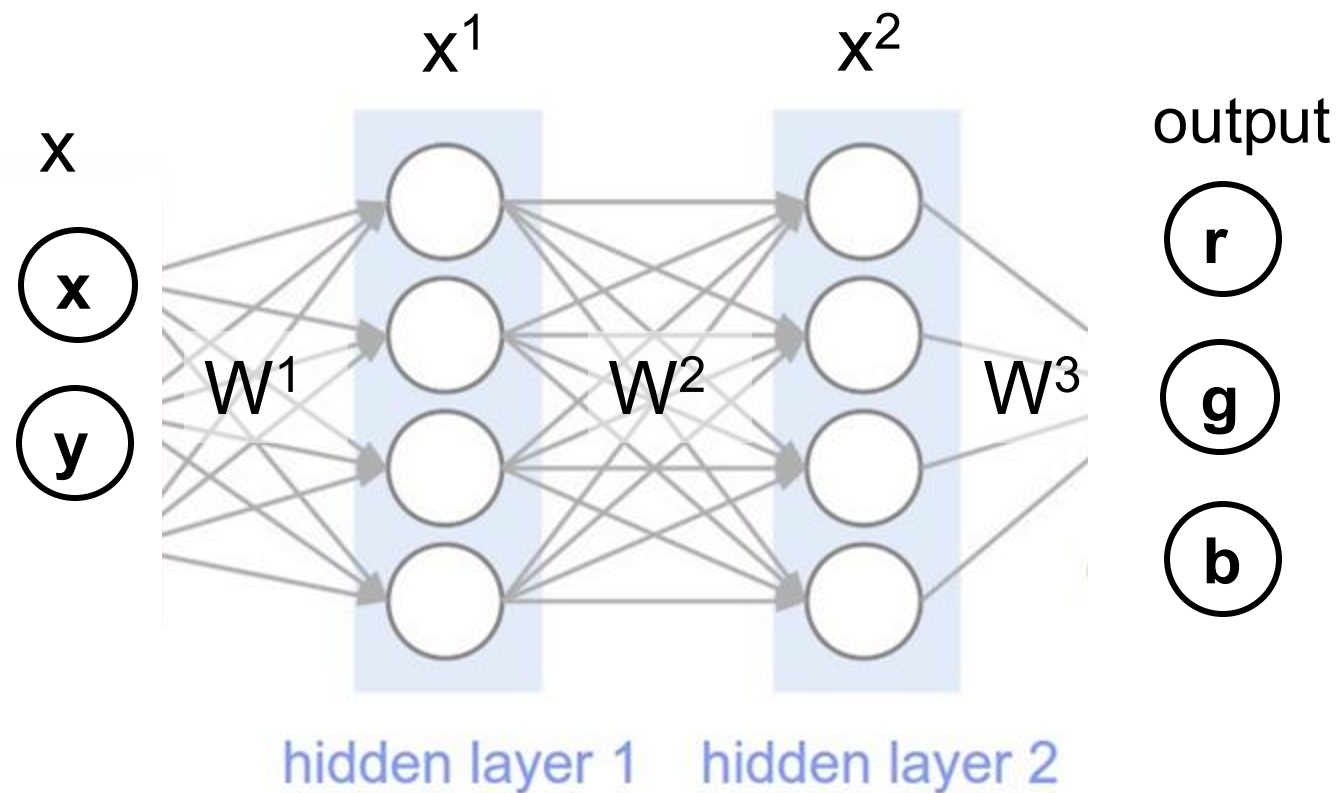
1. Linear Transform $z = W^l x^{l-1} + b$
2. Apply Non-Linearity $x^l = f(z)$

Usually

$$f = \text{RELU}(z) \\ = \max(0, z)$$

What are the learnable parameters?

In our 2D case:



In each layer:

1. Linear Transform $z = W^l x^{l-1} + b$
2. Apply Non-Linearity $x^l = f(z)$

Usually

$$f = \text{RELU}(z) \\ = \max(0, z)$$

What are the learnable parameters?

Coordinate Based Neural Network

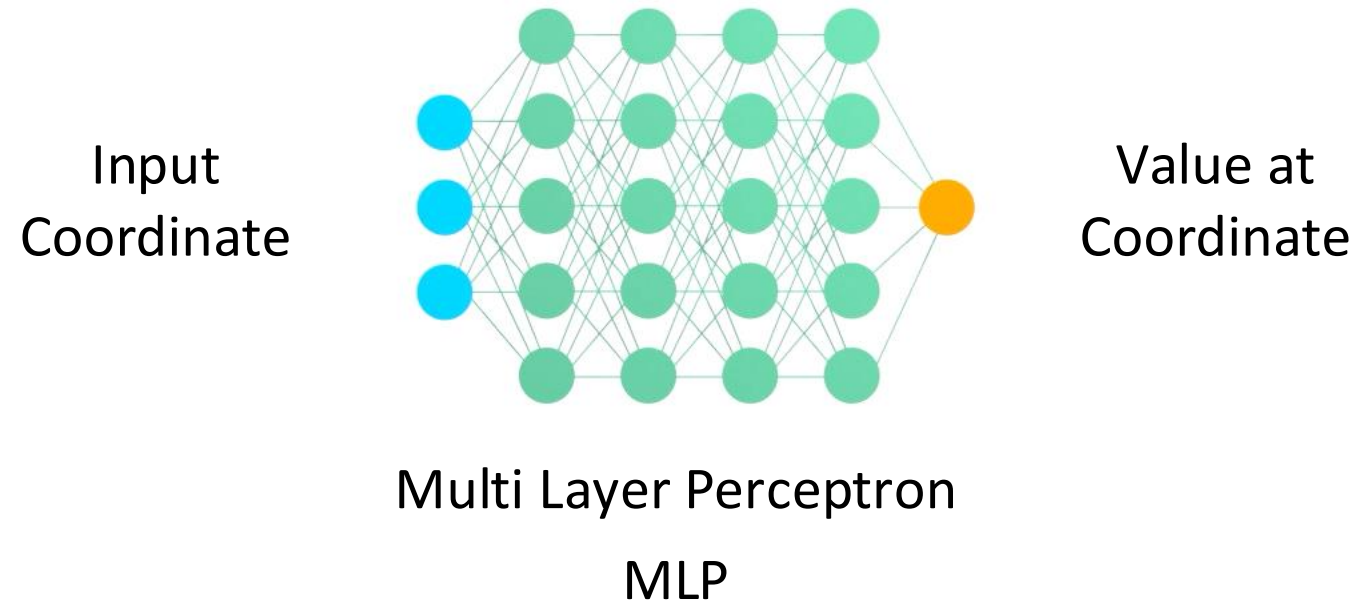
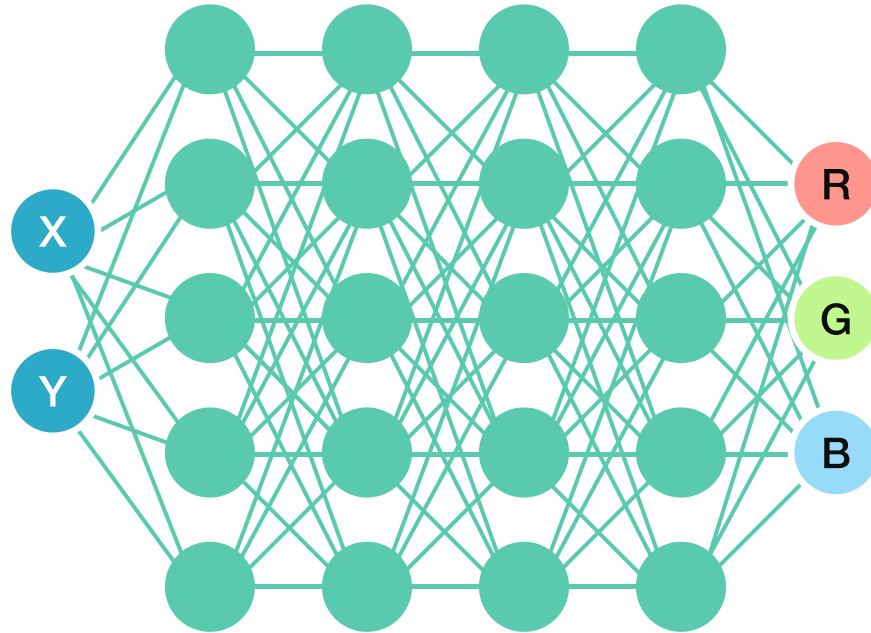


Image Representation



Challenge:

How to get MLPs to represent higher frequency functions?

what happens if you naively
optimize this network

Iteration 1000

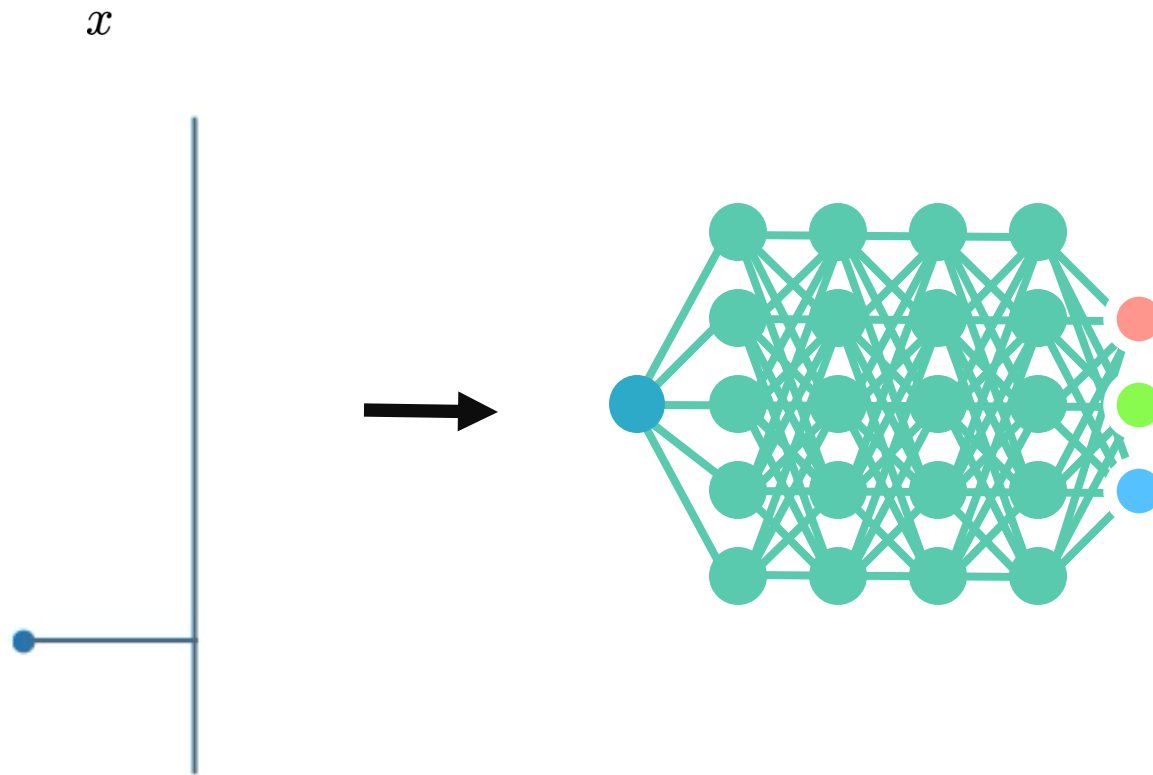


MLP output



Supervision image

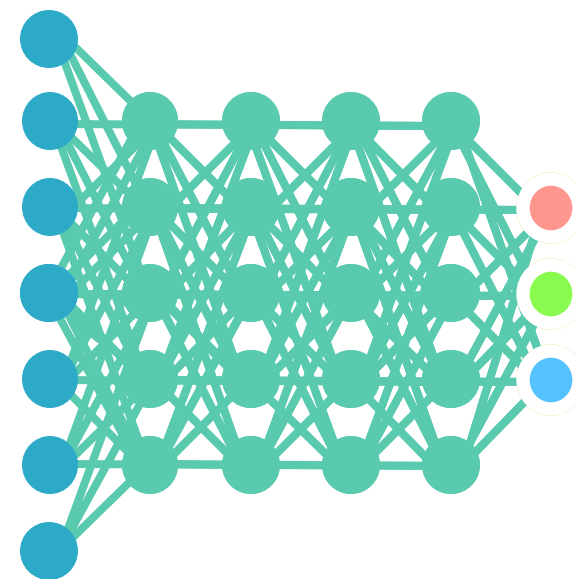
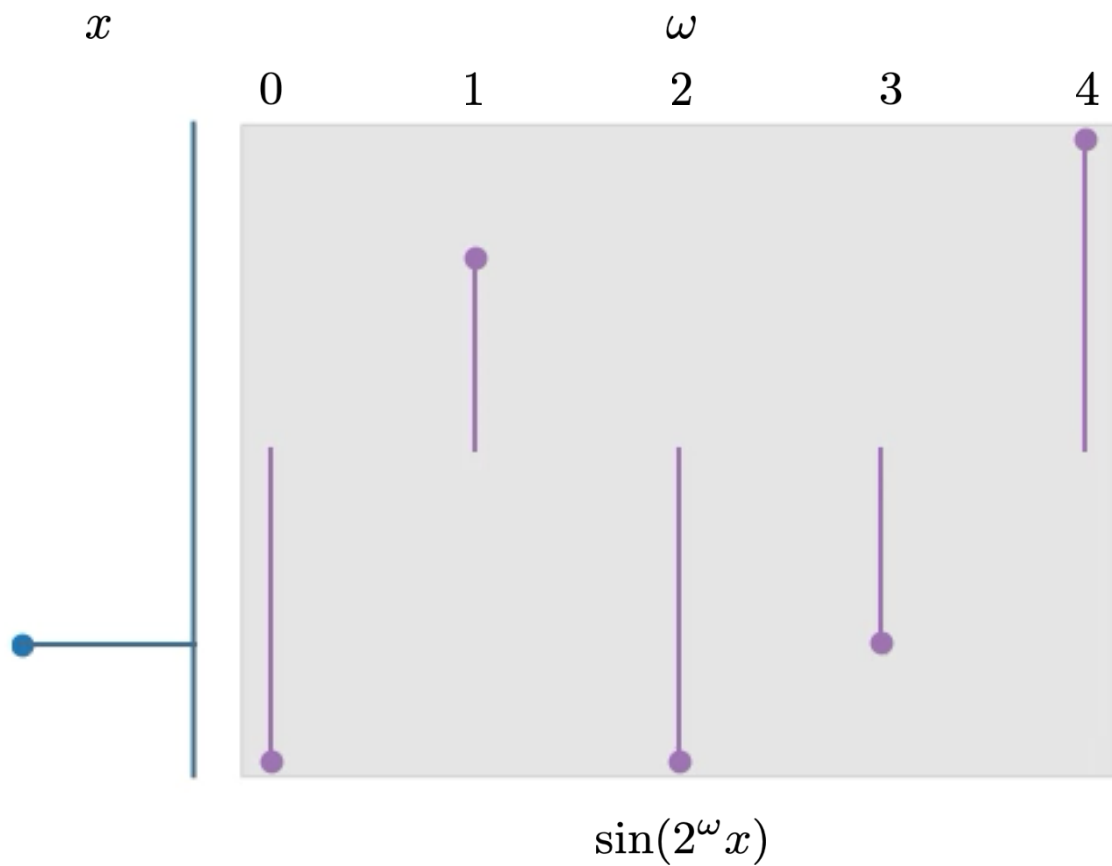
Standard input



Positional Encoding

Standard input

Positionally Encoded input



Fourier Features

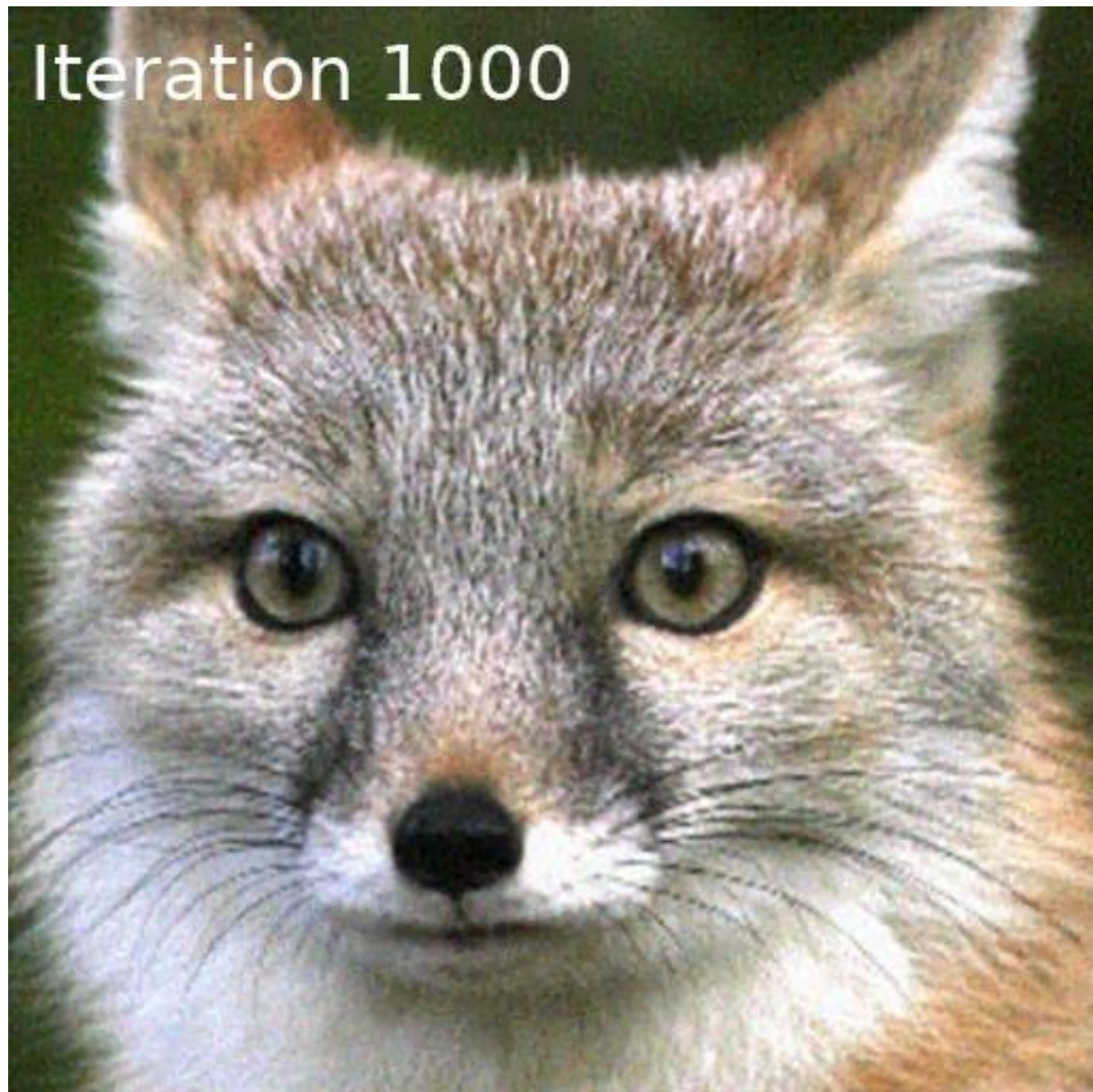
$$\gamma(p) = (\sin(2^0 \pi p), \cos(2^0 \pi p), \dots, \sin(2^{L-1} \pi p), \cos(2^{L-1} \pi p))$$

Iteration 1000



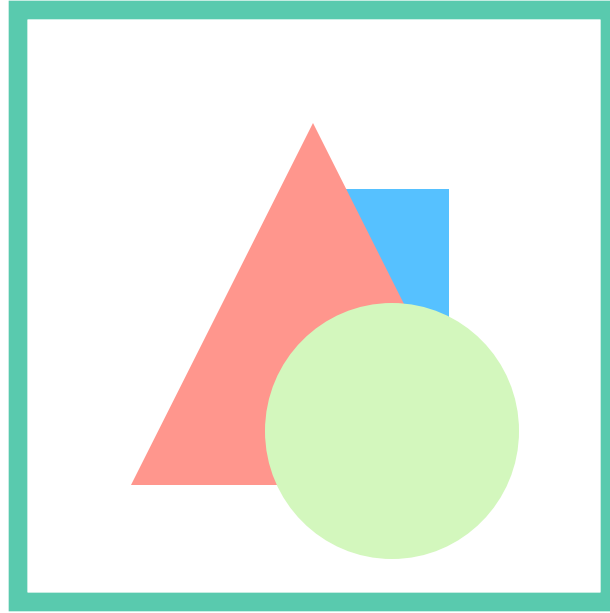
Standard MLP

Iteration 1000



MLP with Fourier features

Why does positional encoding help?



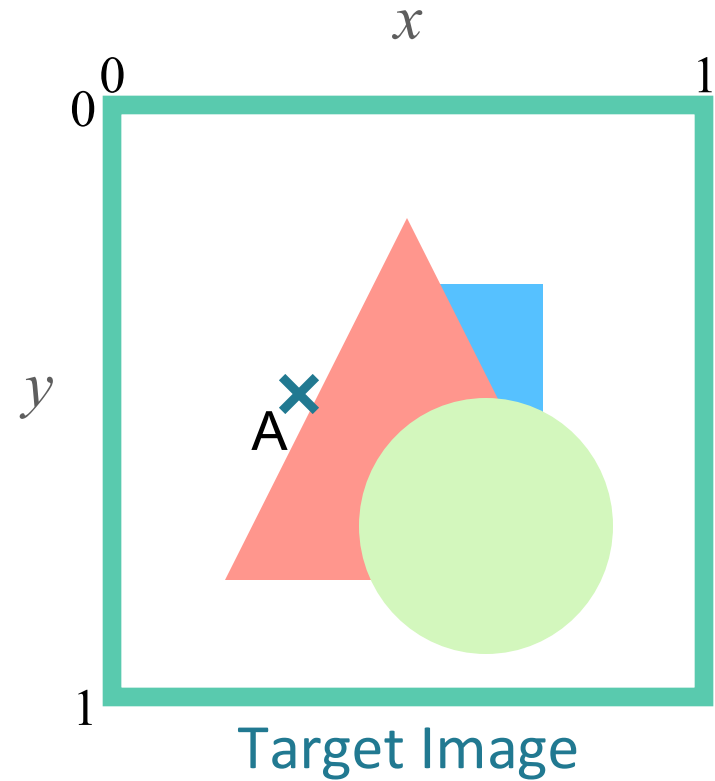
Target Image

Why does positional encoding help?

Input

	x	y
A	.36	.5

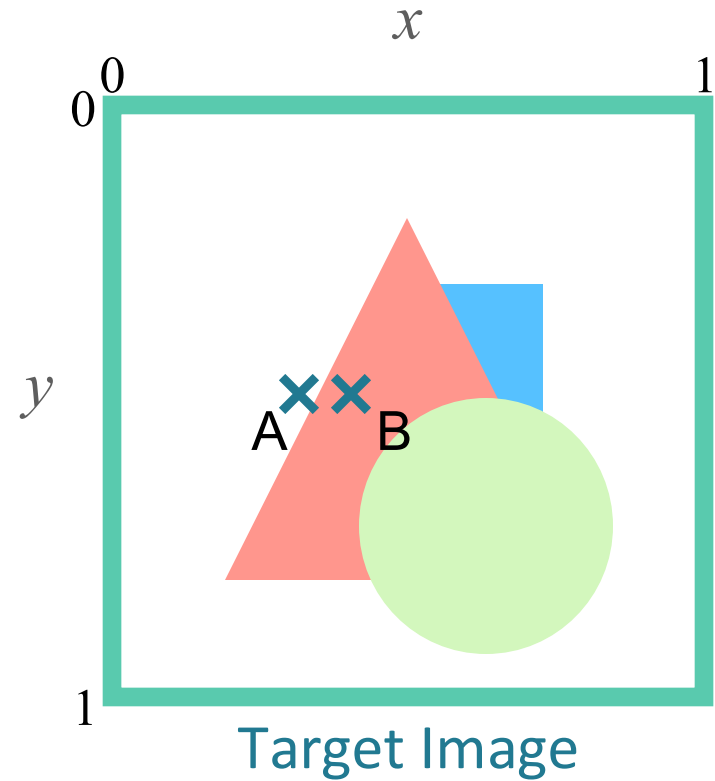
Target



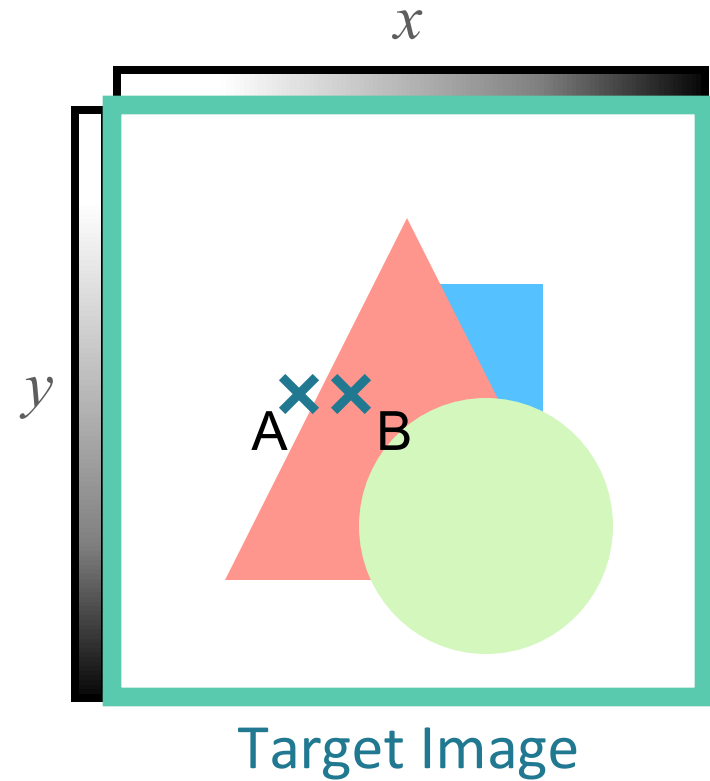
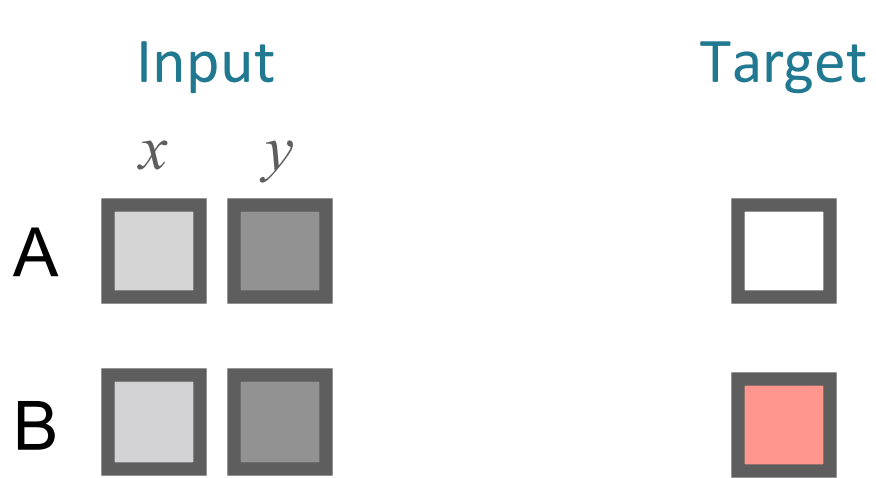
Why does positional encoding help?

Input		
	x	y
A	.36	.5
B	.38	.5

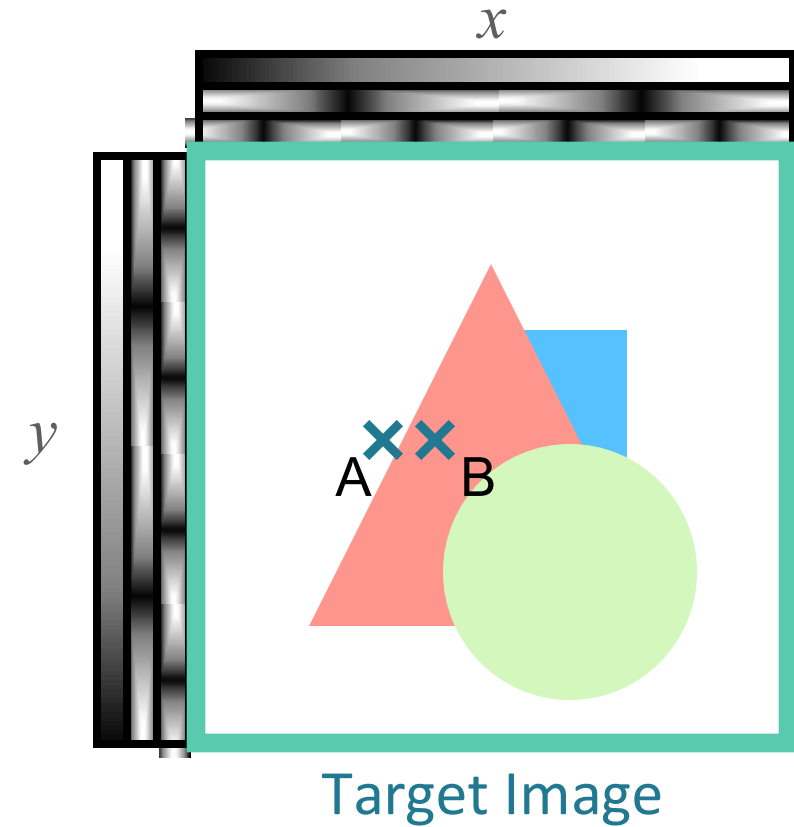
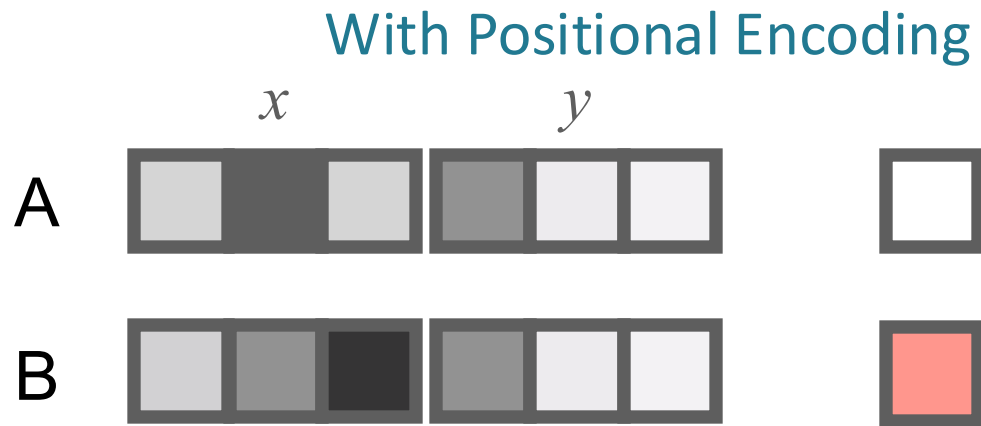
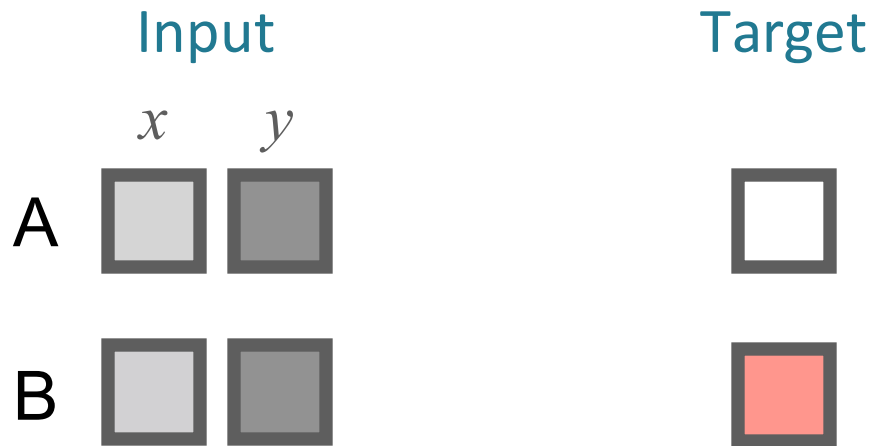
Target



Why does positional encoding help?

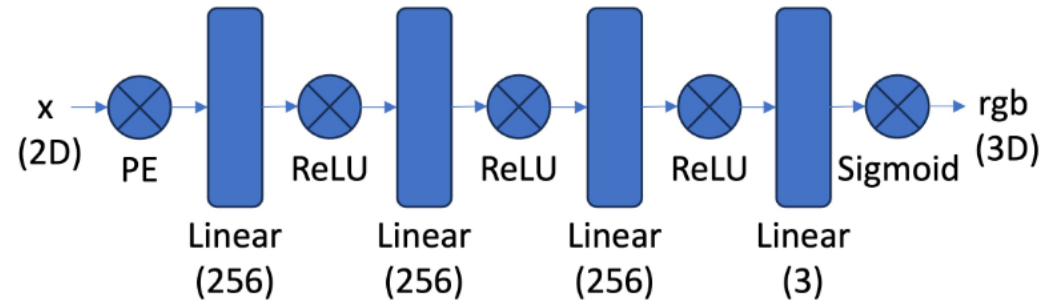


Why does positional encoding help?



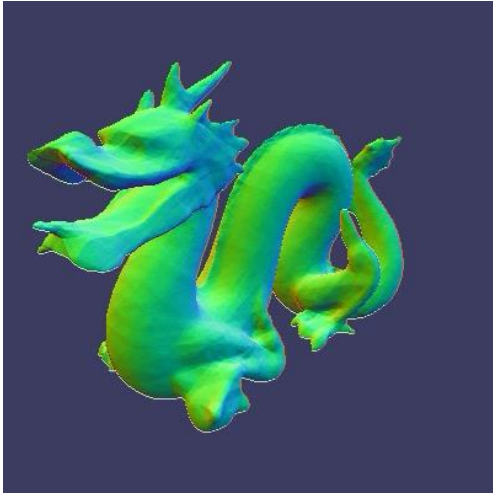
NeRF Project Part 1

- Fit a Neural Network to a single image
- Implement this network, and Positional Embedding (PE) and reconstruct an image:

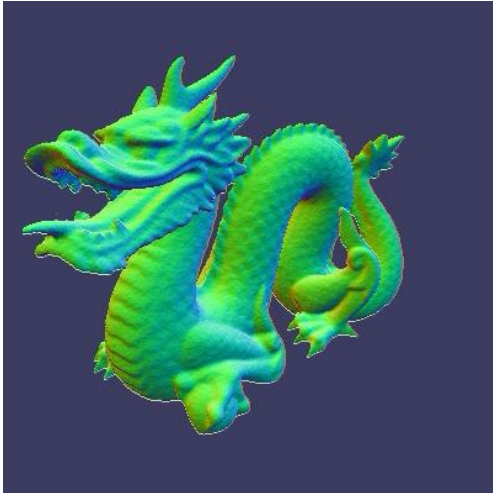


Coordinate-based MLPs can replace any low-dimensional array

Without Encoding



With Encoding



3D Shape

NeRF with and without positional encoding



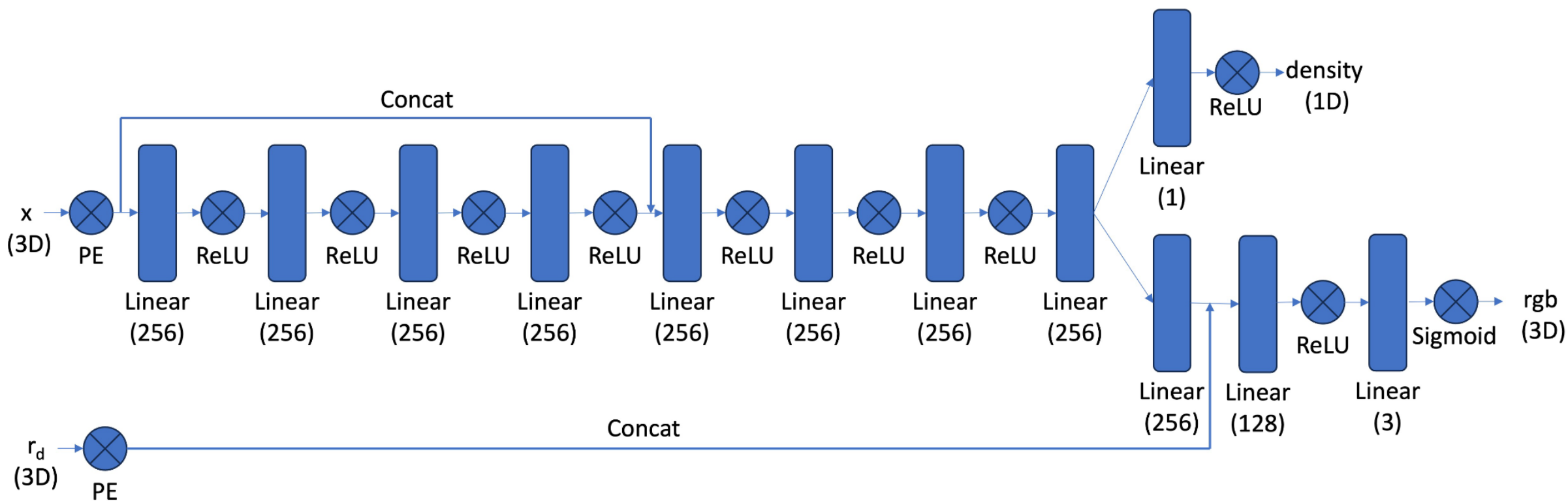
NeRF (Naive)



NeRF (with positional encoding)

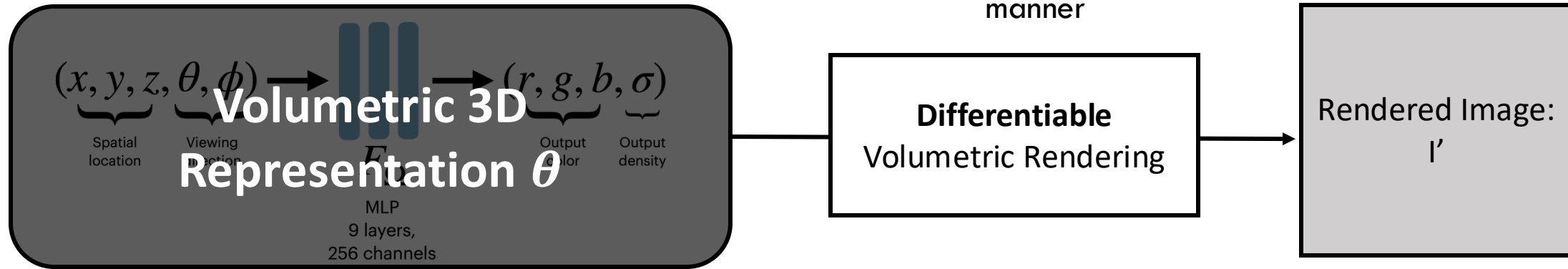
NeRF Network Architecture

Next section you will implement this:



Let's go back to 3D

Now we need to render an image from this 3D representation in a differentiable manner

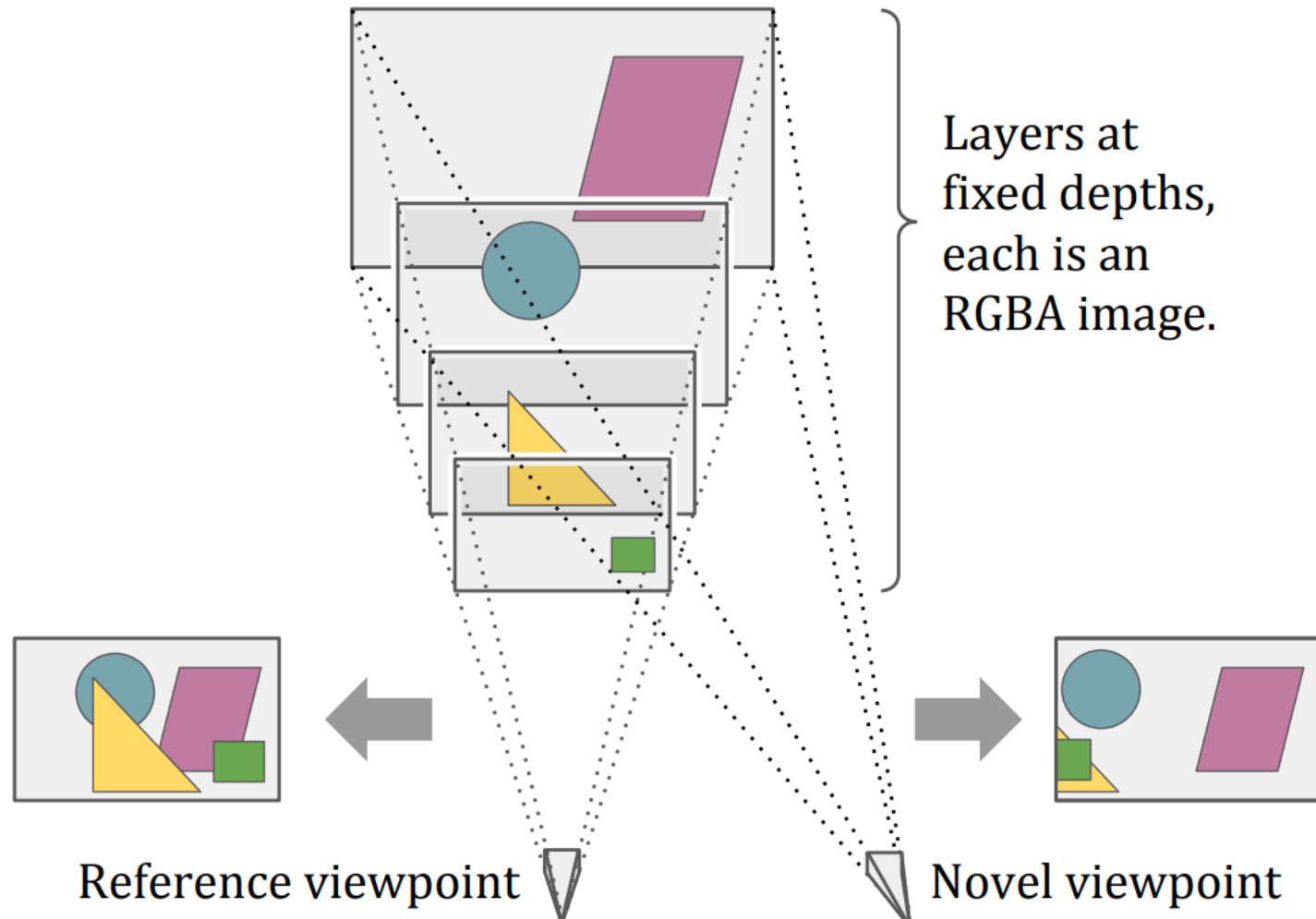


“Training” Objective (aka Analysis-by-Synthesis):

$$\min_{\theta} \| \text{Rendered Image: } I' - \text{Observed Image: } I \|_2$$

Differentiable Volumetric Rendering

A Precursor: Multi-plane Images



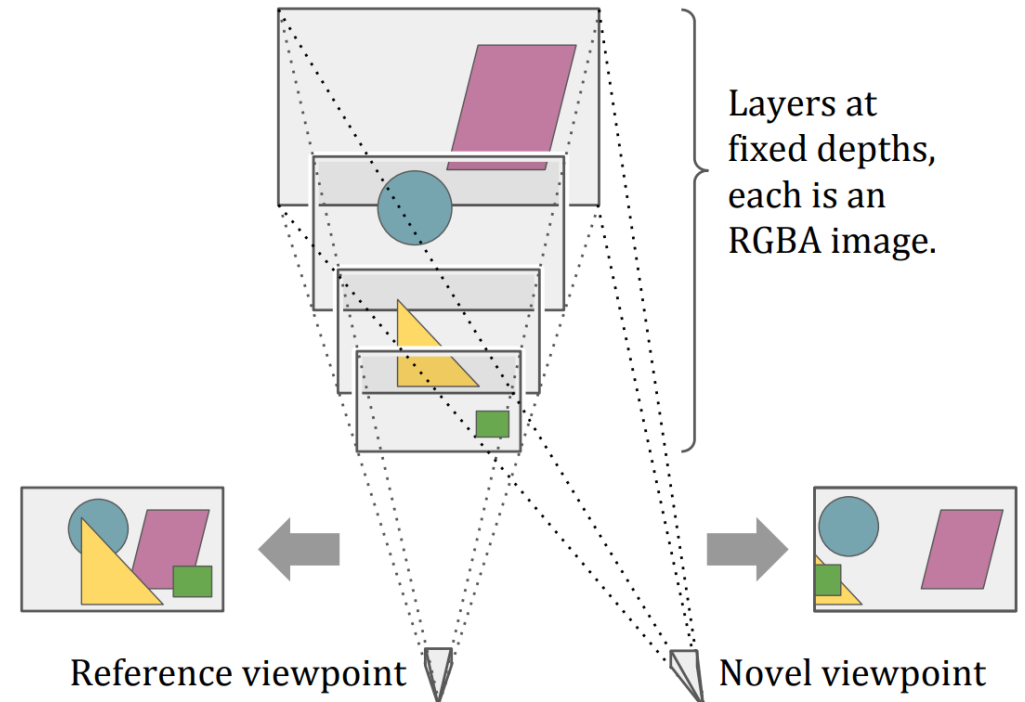
Multi- plane Camera at Disney

<https://www.youtube.com/watch?v=YdHTIUGN1zw>

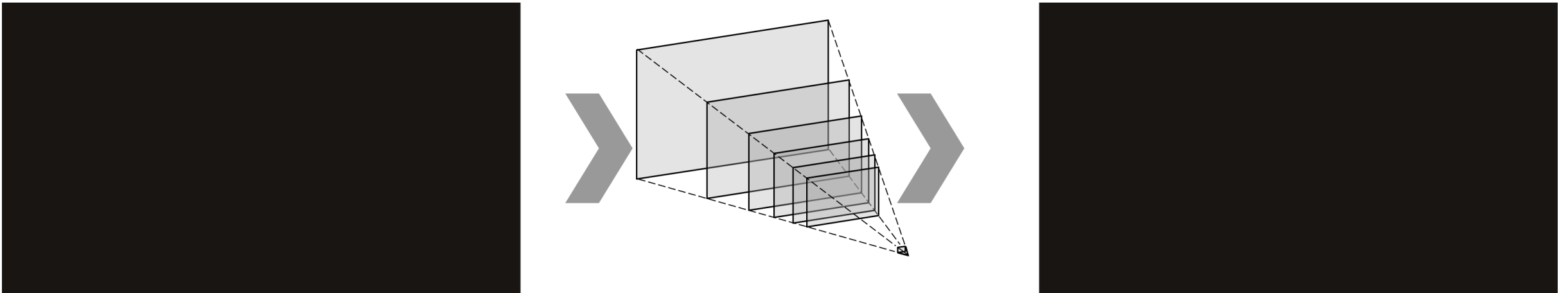
Generating an Image MPI

To render a novel view:

1. Homography warp the image from the new viewpoint
2. Alpha Blend each layer



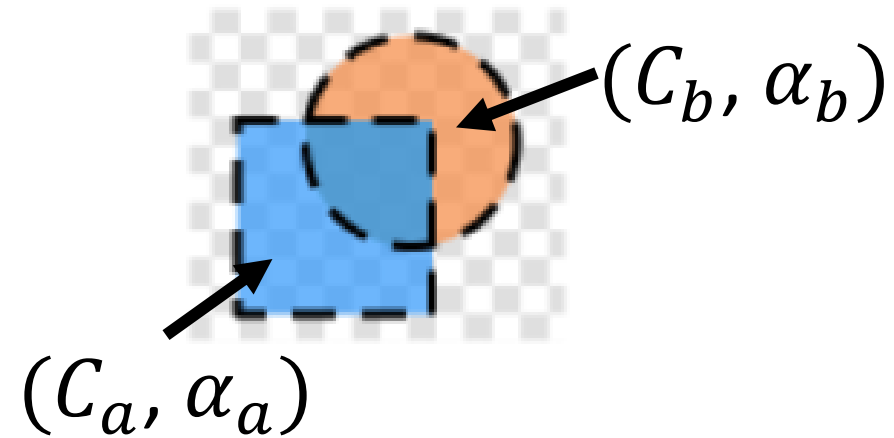
Sample Novel View Synthesis with a MPI



Also called front-to-back compositing or “over” operation

Alpha Blending

for two image case, A and B, both partially transparent:

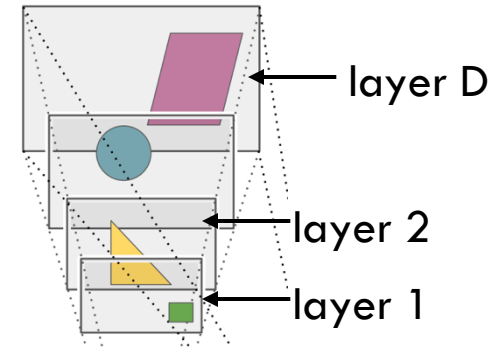


$$I = C_a \alpha_a + C_b \alpha_b \underbrace{(1 - \alpha_a)}$$

How much light is the previous layer letting through?

General D layer case:

$$I = \sum_{i=1}^D C_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j)$$



What is missing in MPIs?

- Look at it from the side??
- You'll see all the edges!!

➔ Limited camera mobility

NeRF overcomes this problem, because it's defined everywhere
Volumetric Rendering behaves similarly to alpha compositing

Back to NeRFs

Neural Volumetric Rendering

Through Volumetric
Representation
(No surfaces)!

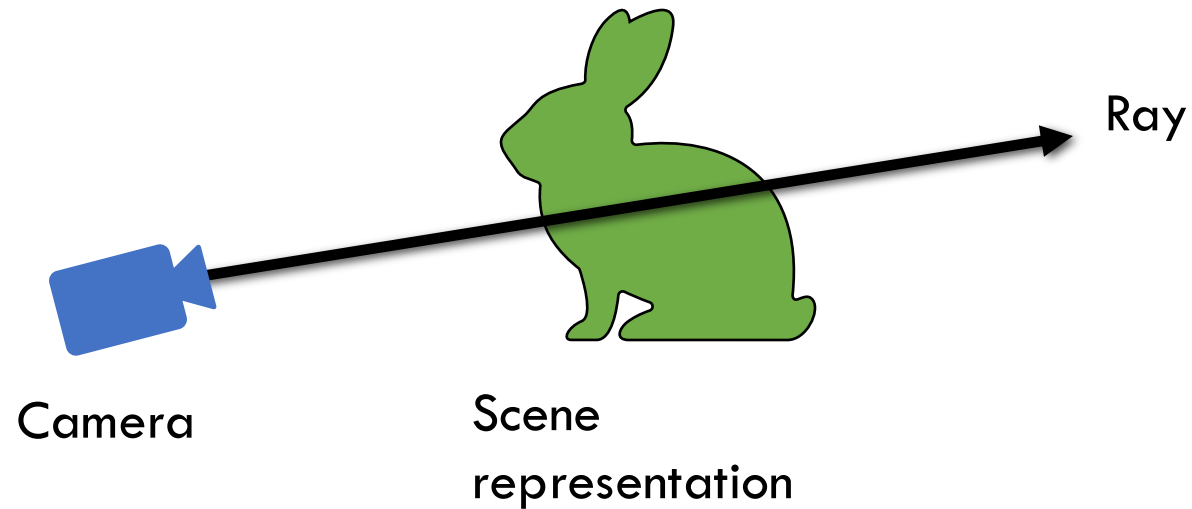


computing color along rays
through 3D space



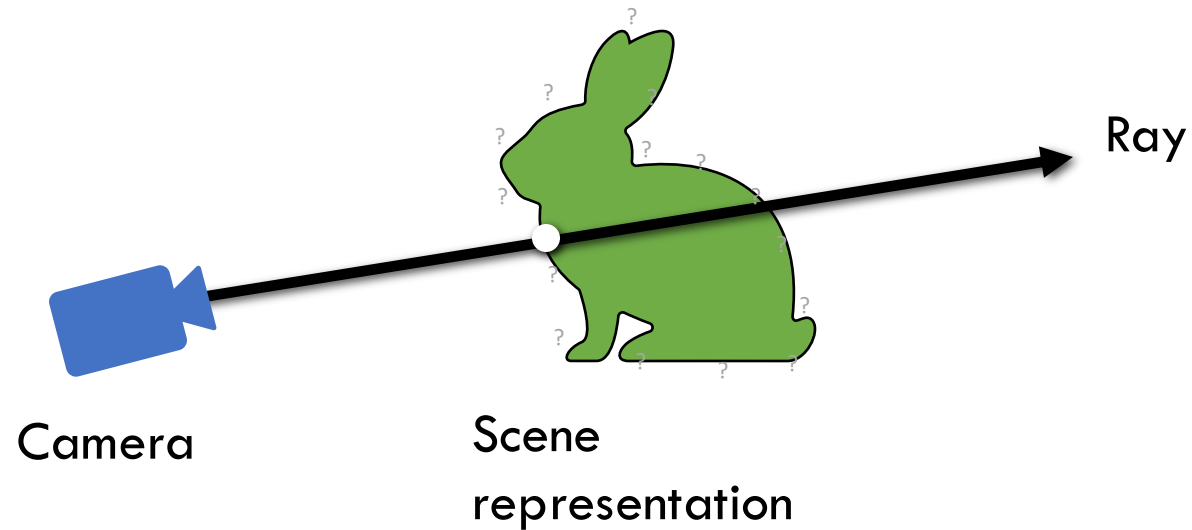
What color is this pixel?

Surface vs. volume rendering



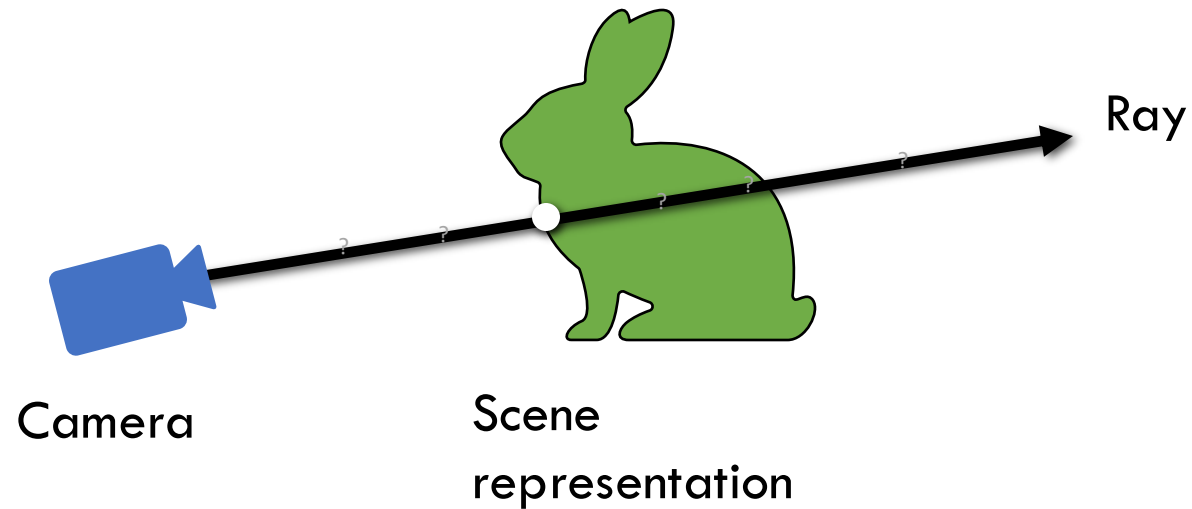
Want to know how ray interacts with scene

Surface vs. volume rendering



Surface rendering — loop over geometry, check for ray hits

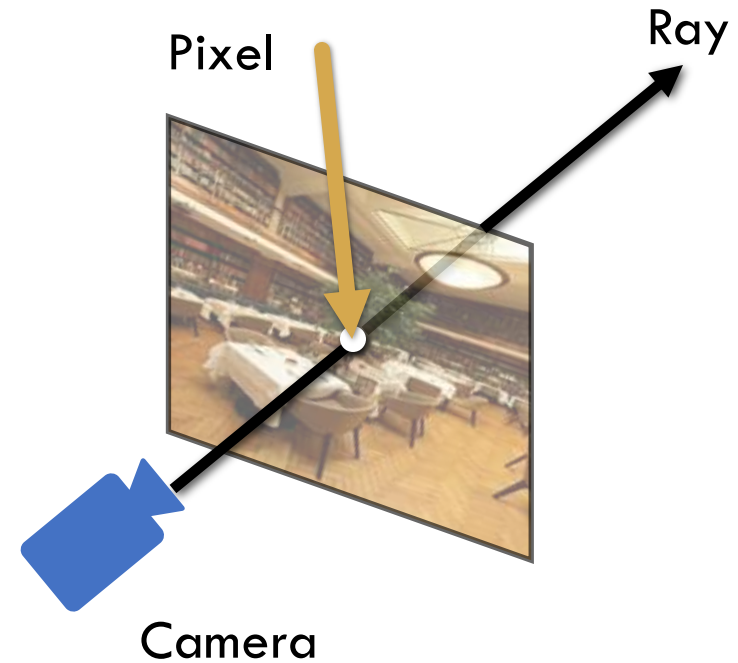
Surface vs. volume rendering



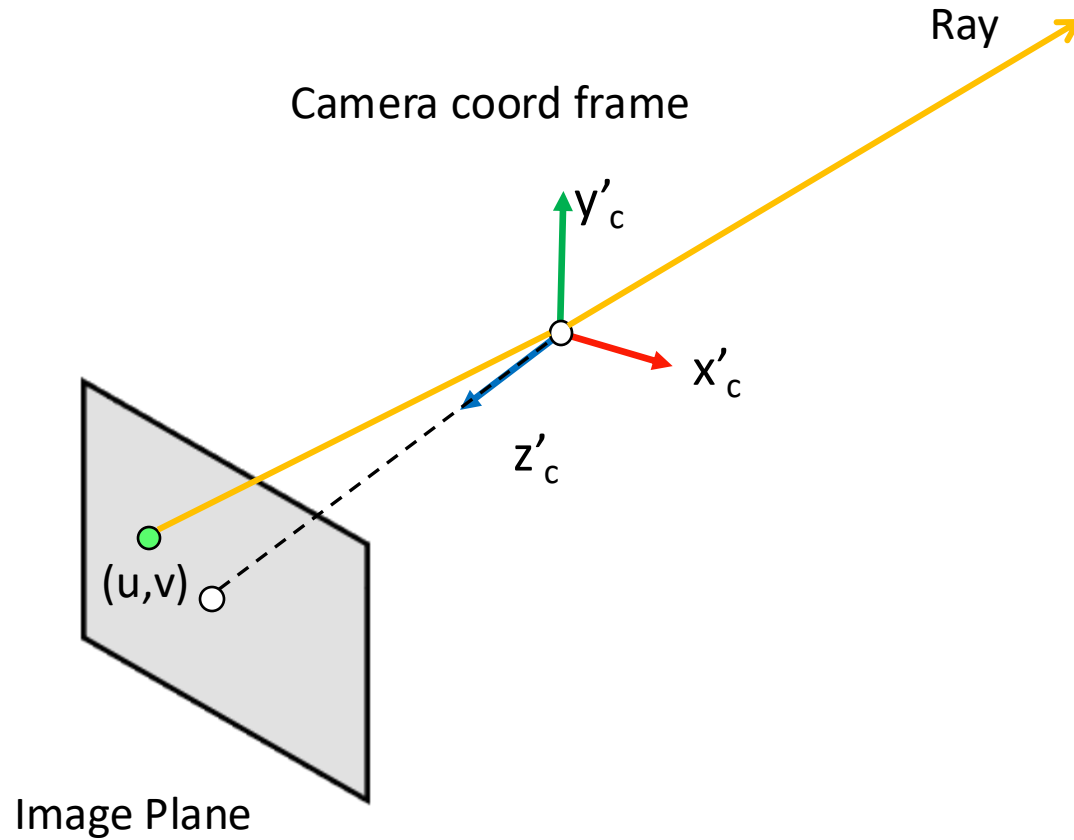
Volume rendering — loop over ray points, query geometry

Recap: Cameras and rays

- We need the mathematical mapping from $(camera, pixel) \rightarrow ray$
- Then can abstract underlying problem as learning the function $ray \rightarrow color$



Compute the Ray



3D to 2D:
(point)

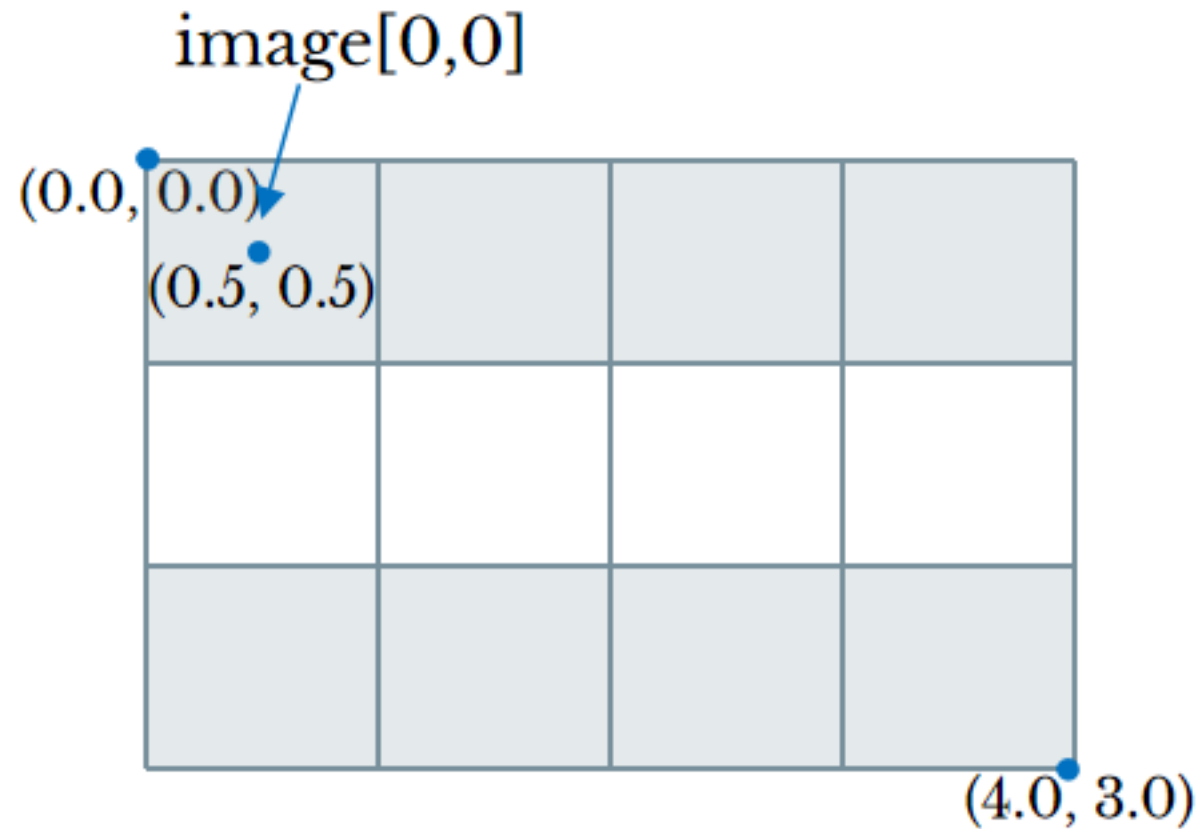
$$u = f_x \frac{x_c}{z_c} + o_x$$
$$v = f_y \frac{y_c}{z_c} + o_y$$

2D to 3D:
(ray)
Back projection

$$x = \frac{z}{f_x} (u - o_x)$$
$$y = \frac{z}{f_y} (v - o_y)$$
$$z > 0$$

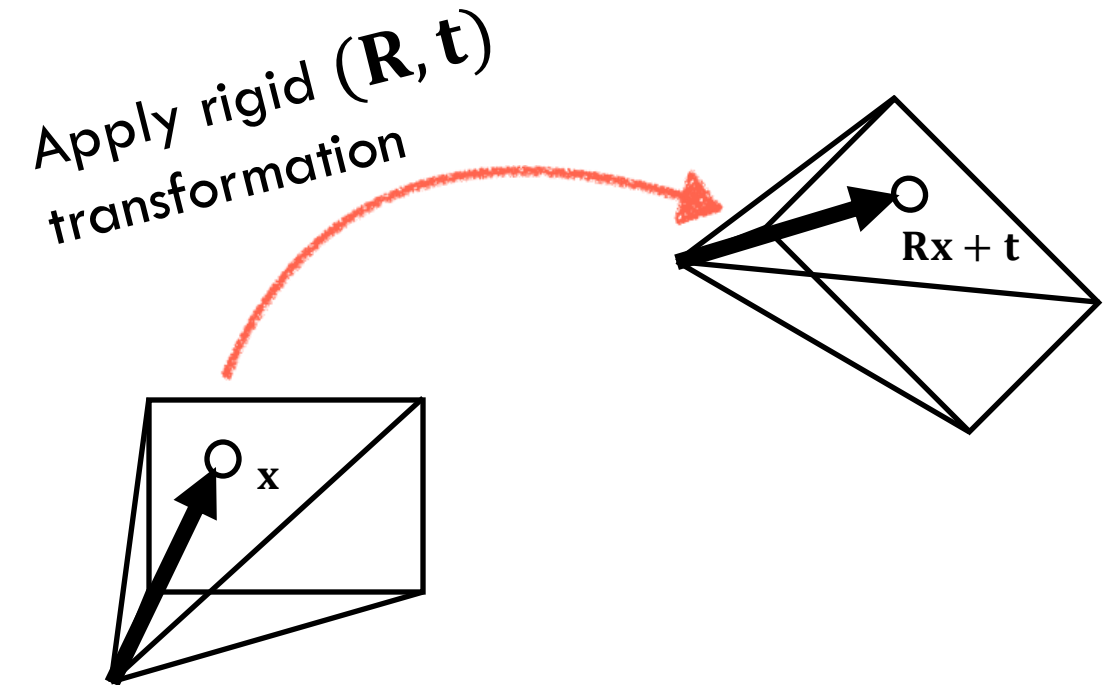
Details:

A half-pixel offset — add 0.5 to i and j so ray precisely hits pixel center



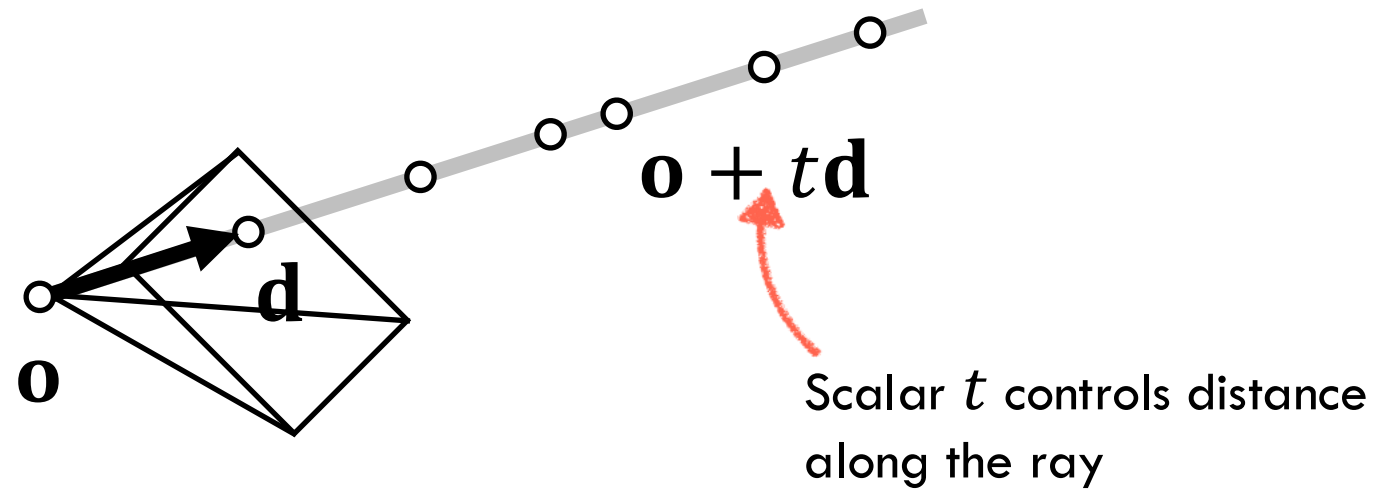
Want: Ray in the World

- What coordinate space is the current ray in?
- Convert it to World!



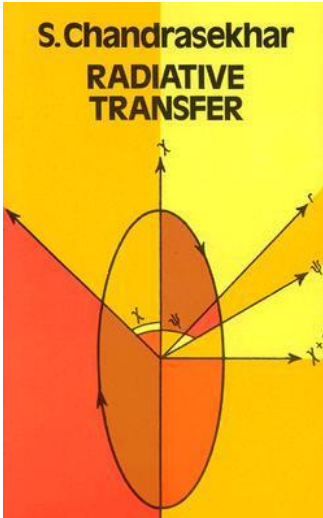
Calculating points along a ray

In the world coordinate frame:

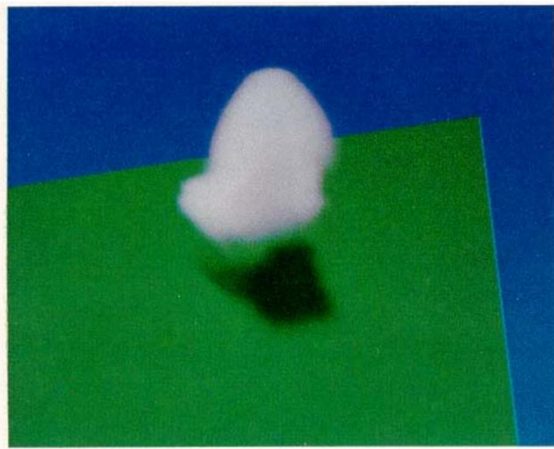


History of volume rendering

In Early computer graphics

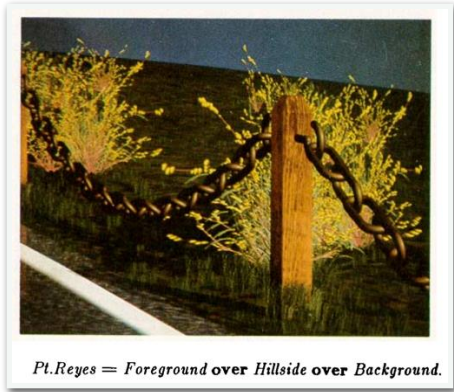


- ▶ Theory of volume rendering co-opted from physics in the 1980s: absorption, emission, out-scattering/in-scattering
- ▶ Adapted for visualising medical data and linked with alpha compositing
- ▶ Modern path tracers use sophisticated Monte Carlo methods to render volumetric effects



Ray tracing simulated cumulus cloud [Kajiya]

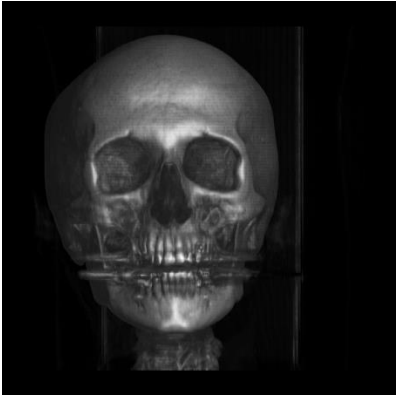
Alpha compositing



Alpha compositing [Porter and Duff]

- ▶ Theory of volume rendering co-opted from physics in the 1980s: absorption, emission, out-scattering/in-scattering
- ▶ Alpha rendering developed for digital compositing in VFX movie production

Volume rendering for visualization



Medical data visualisation [Levoy]

- ▶ Theory of volume rendering co-opted from physics in the 1980s: absorption, emission, out-scattering/in-scattering
- ▶ Alpha rendering developed for digital compositing in VFX movie production
- ▶ **Volume rendering applied to visualise 3D medical scan data in 1990s**

Chandrasekhar 1950, Radiative Transfer

Kajiya 1984, Ray Tracing Volume Densities

Porter and Duff 1984, Compositing Digital Images

Levoy 1988, Display of Surfaces from Volume Data

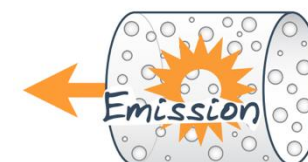
Max 1995, Optical Models for Direct Volume Rendering



Absorption



Scattering



Emission



<http://commons.wikimedia.org>



<http://wikipedia.org>

Simplify

Absorption



<http://commons.wikimedia.org>

Scattering



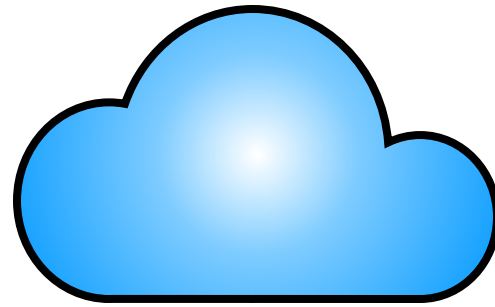
Emission



<http://wikipedia.org>

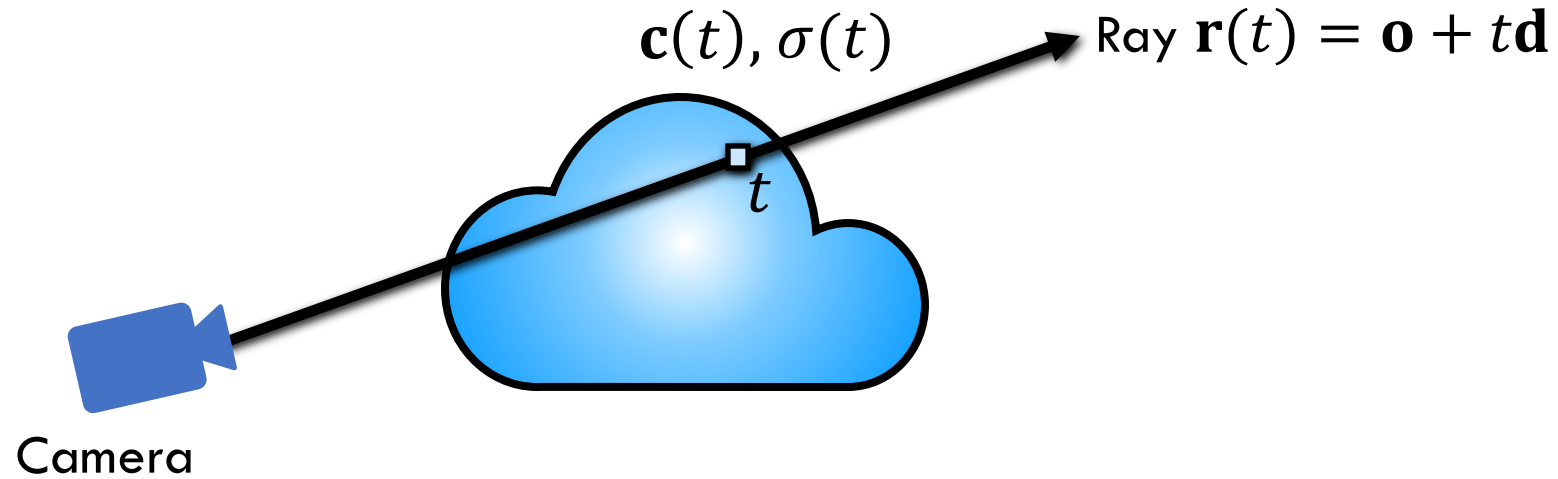
Volume rendering derivations

Volumetric formulation for NeRF



Scene is a cloud of tiny colored particles

Volumetric formulation for NeRF



at a point on the ray $\mathbf{r}(t)$, we can query color $\mathbf{c}(t)$ and density $\sigma(t)$

How to integrate all the info along the ray to get a color per ray?

Idea: Expected Color

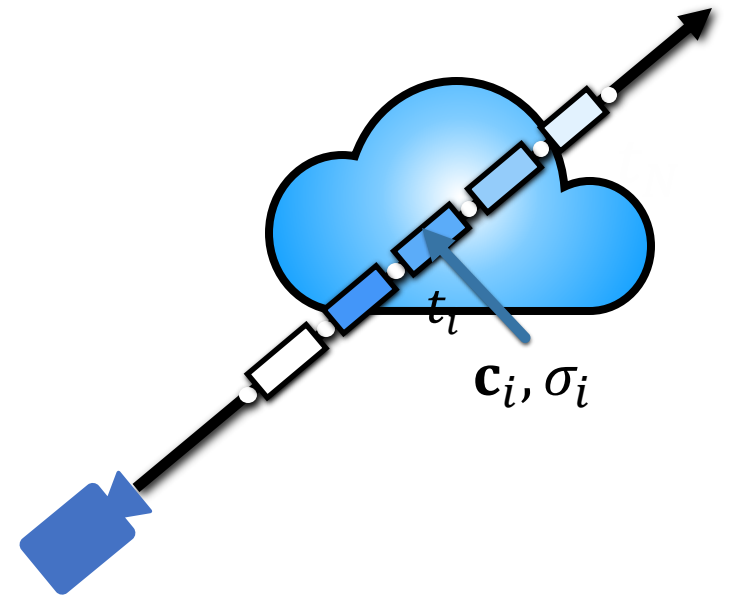
- Pose probabilistically.
- Each point on the ray has a probability to be the first “hit” : $P[\text{first hit at } t]$
- Color per ray = Expected value of color with this probability of first “hit”

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

$$\mathbf{c}(\mathbf{r}) = \int_{t_0}^{t_1} P[\text{first hit at } t] \mathbf{c}(t) dt$$

$$\approx \sum_{t=0}^T P[\text{first hit at } t] \mathbf{c}(t)$$

$$\approx \sum_{t=0}^T w_t \mathbf{c}(t)$$



Differentiable Volumetric Rendering Formula

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

$$\mathbf{c} \approx \sum_{i=1}^n w_i \mathbf{c}_i$$

differentiable w.r.t. \mathbf{c}, σ

weights

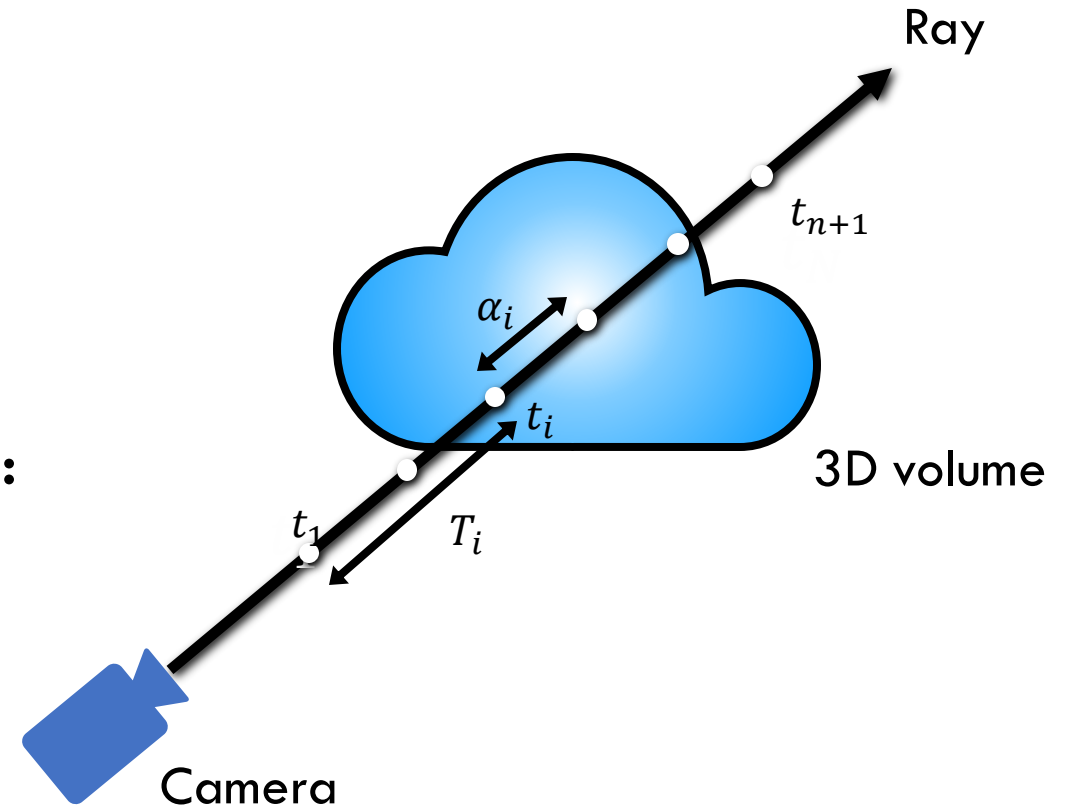
colors

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 :

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



Summary

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

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

weights colors

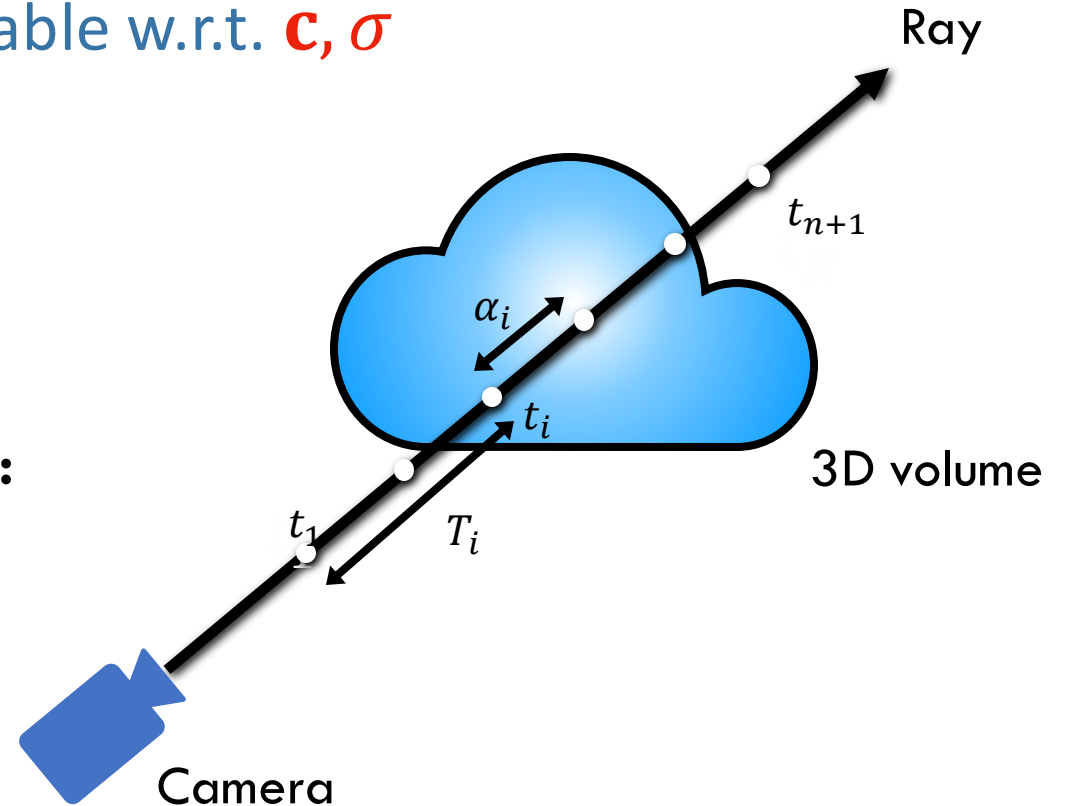
differentiable w.r.t. \mathbf{c}, σ

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 :

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

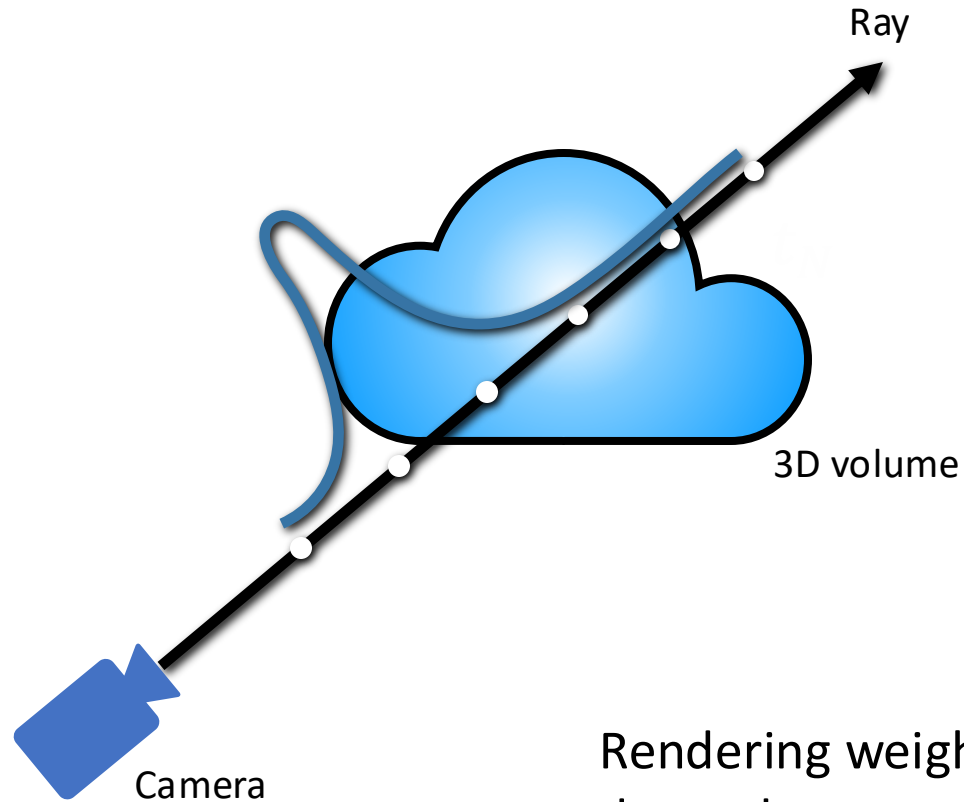


Complete derivation

- If you want a complete derivation of the volrend equation, see this <https://drive.google.com/file/d/1QsCK5V0d6DSc0QGcKsV97u83JFd-dFoz/view>
- From slide 35
- Or <https://arxiv.org/pdf/2209.02417>

Visual intuition: rendering weights is specific to a ray

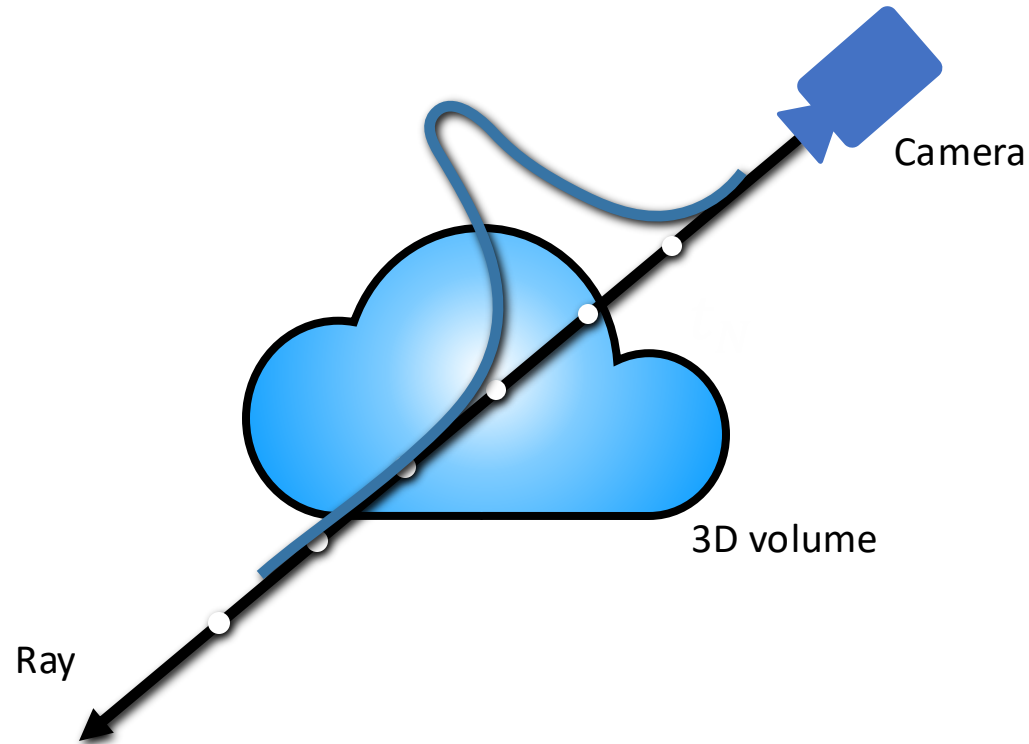
$$C \approx \sum_{i=1}^N T_i \alpha_i c_i$$



Rendering weights are not a 3D function — depends on ray, because of transmittance!

Visual intuition: rendering weights is specific to a ray

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



Rendering weights are not a 3D function — depends on ray, because of transmittance!

Rendering weight PDF is important

Remember, expected color is equal to

$$\int T(t)\sigma(t)\mathbf{c}(t)dt \approx \sum_i T_i\alpha_i\mathbf{c}_i = \sum_i w_i\mathbf{c}_i$$

$T(t)\sigma(t)$ and $T_i\alpha_i$ are “rendering weights” — probability distribution along the ray
(continuous and discrete, respectively)

You can also render entities other than color in 3D, for example it's depth, or any other N-D vector \mathbf{v}_i

$$\text{Volume rendered "feature"} = \sum_i w_i \mathbf{v}_i$$

Rendering weight PDF is important — depth

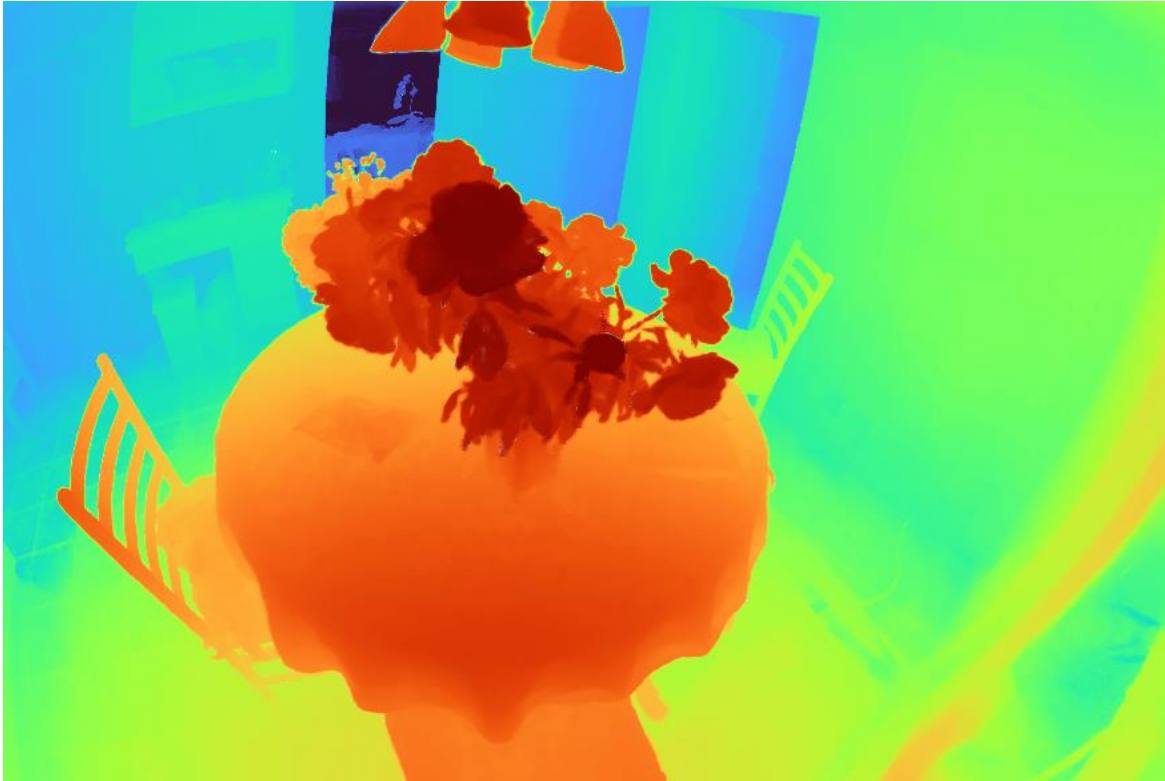
We can use this distribution to compute expectations for other quantities, e.g. “expected depth”:

$$\bar{t} = \sum_i T_i \alpha_i t_i$$

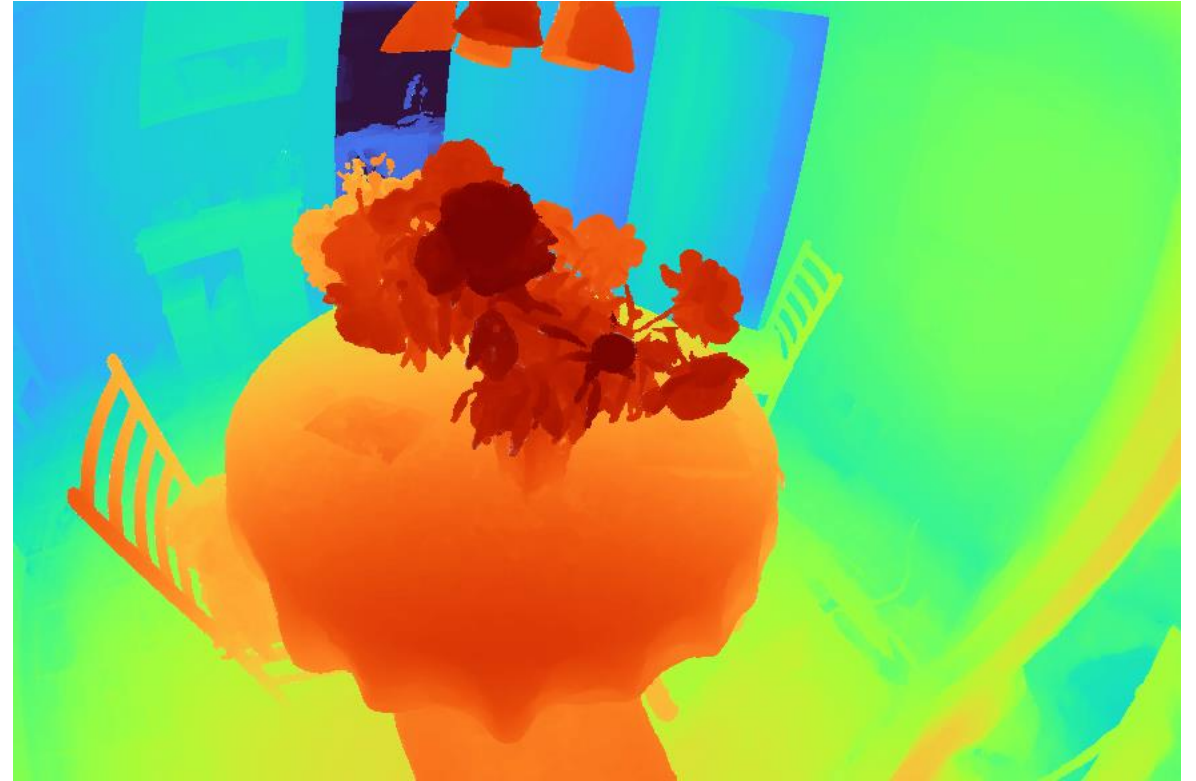
This is often how people visualise NeRF depth maps.

Alternatively, other statistics like mode or median can be used.

Rendering weight PDF is important — depth



Mean depth



Median depth

Rendering weight PDF is important — depth



Mean depth



Median depth



Volume rendering other quantities

This idea can be used for any quantity we want to “volume render” into a 2D image. If \mathbf{V} lives in 3D space (semantic features, normal vectors, etc.)

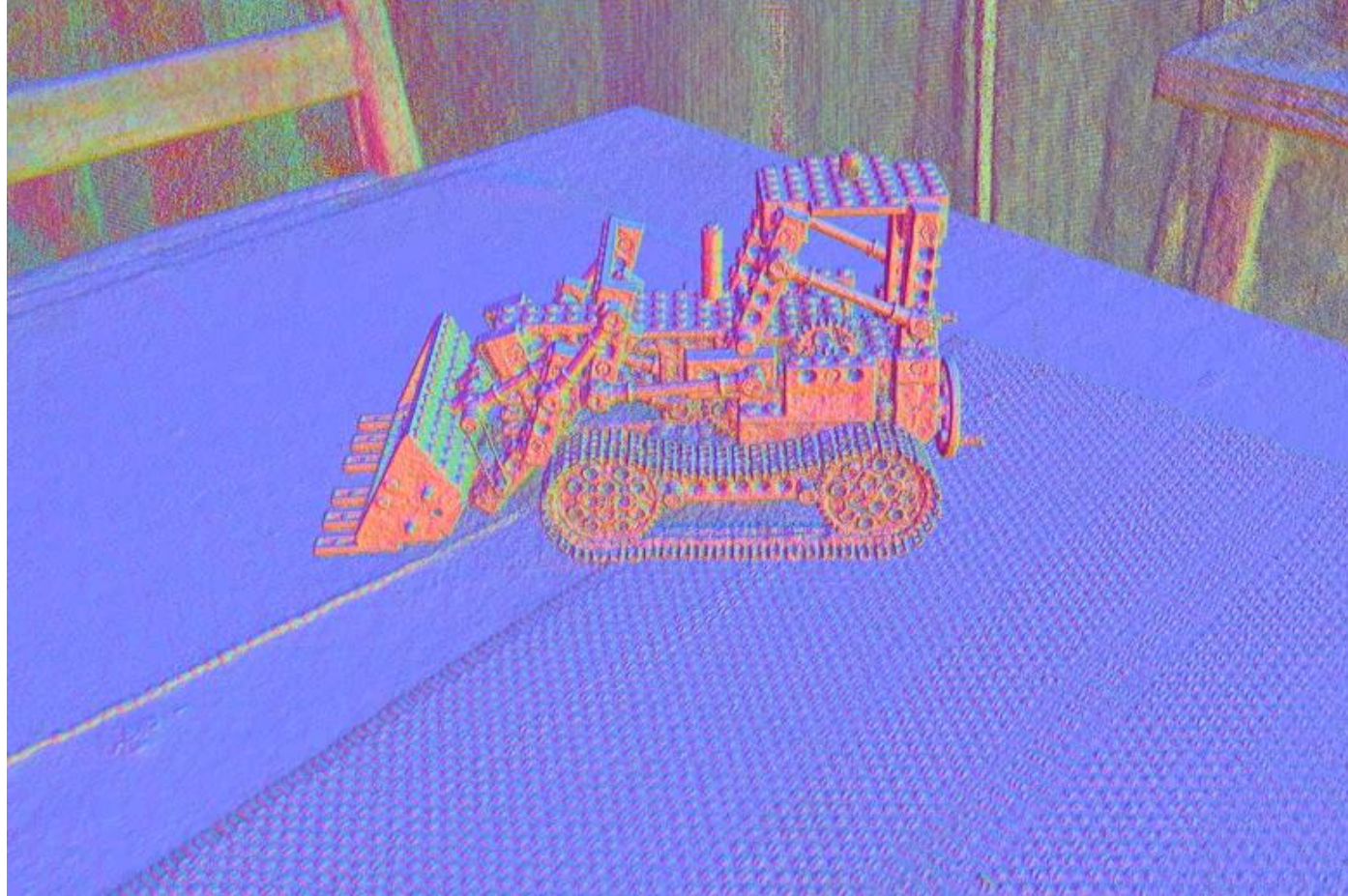
$$\sum_i T_i \alpha_i \mathbf{v}_i$$

can be taken per-ray to produce 2D output images.

Volume Rendering CLIP features

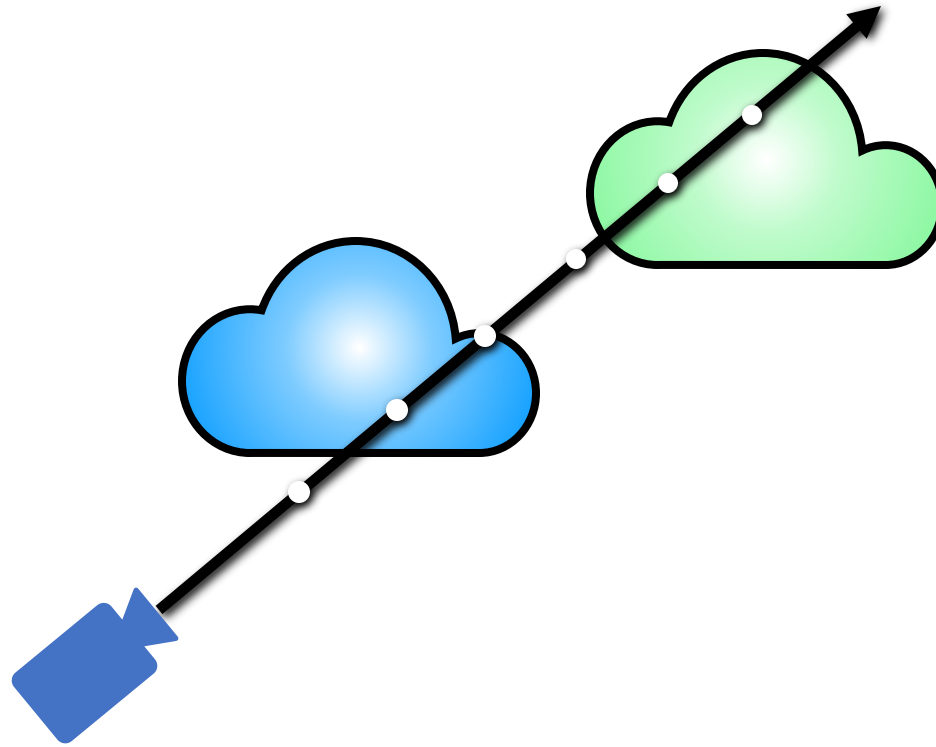


Density as geometry

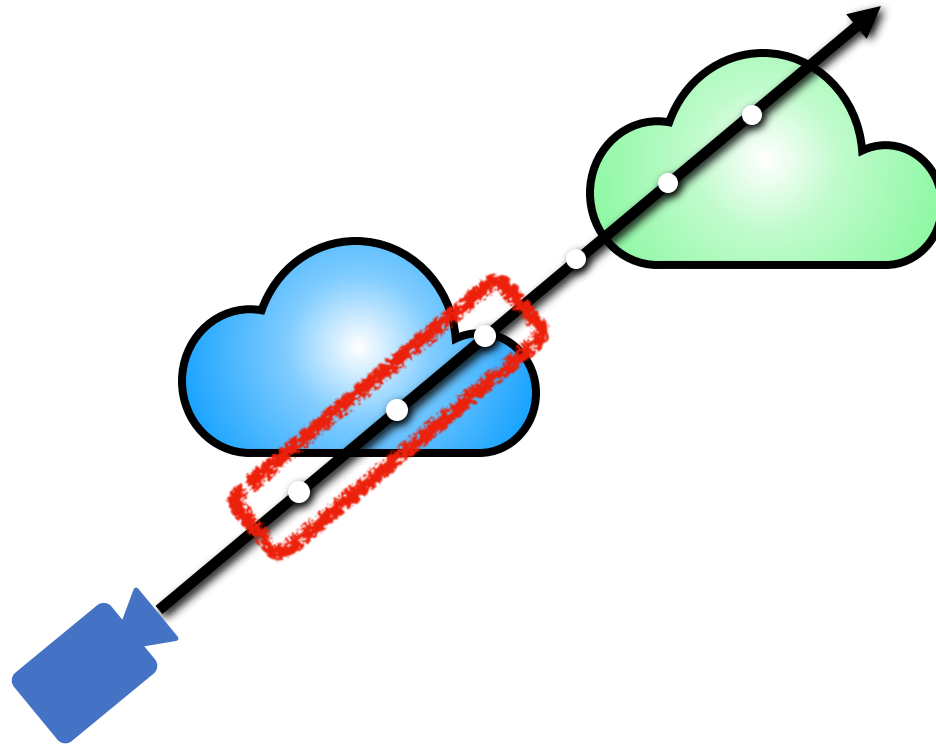


Normal vectors (from analytic gradient of density)

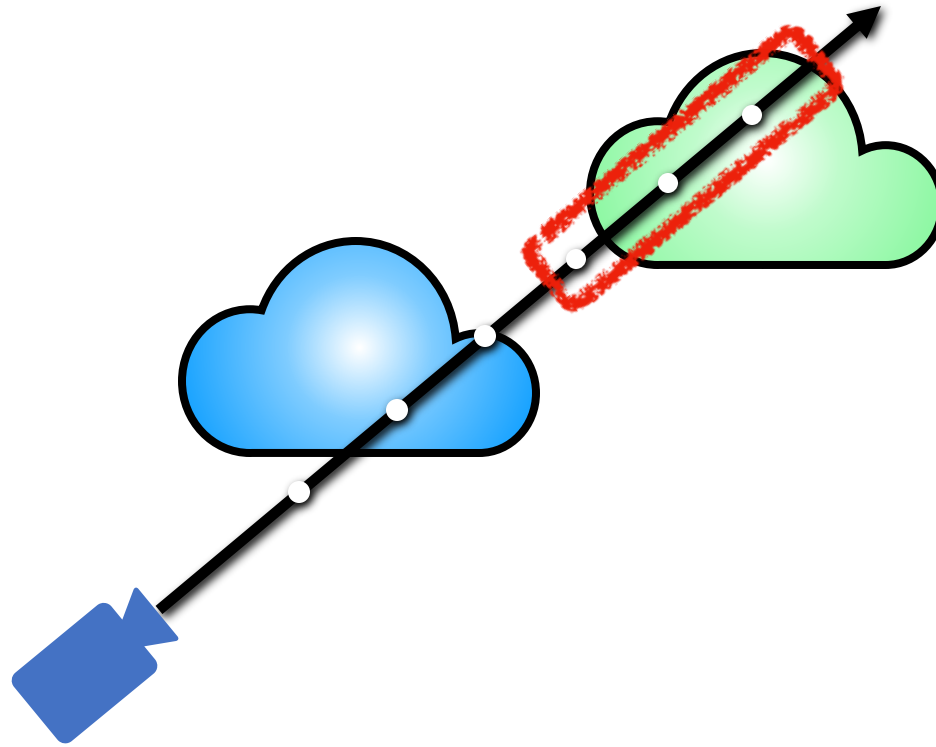
Alpha mattes and compositing



Alpha mattes and compositing



Alpha mattes and compositing



Alpha mattes and compositing

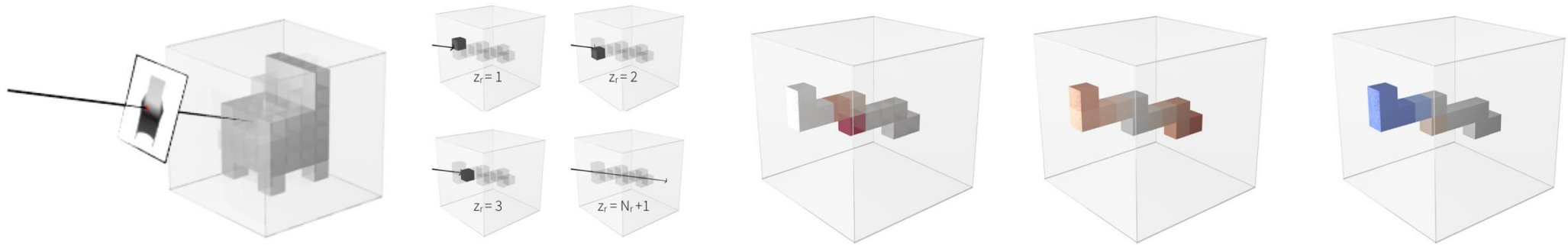


Mildenhall*, Srinivasan*, Tancik* et al 2020, NeRF

Poole et al 2022, DreamFusion

Tang et al 2022, Compressible-composable NeRF via Rank-residual Decomposition

Previous Papers



Differentiable ray consistency work used a forward model with “probabilistic occupancy” to supervise 3D-from-single-image prediction.
Same rendering model as alpha compositing!

$$p(z_r = i) = \begin{cases} (1 - x_i^r) \prod_{j=1}^{i-1} x_j^r, & \text{if } i \leq N_r \\ \prod_{j=1}^{N_r} x_j^r, & \text{if } i = N_r + 1 \end{cases}$$

Similar Ideas before NeRF

Multiplane image methods

Stereo Magnification (Zhou et al. 2018)

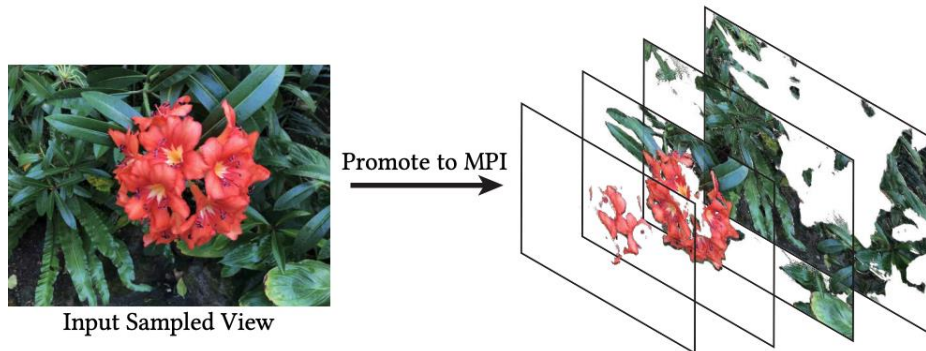
Pushing the Boundaries... (Srinivasan et al. 2019)

Local Light Field Fusion (Mildenhall et al. 2019)

DeepView (Flynn et al. 2019)

Single-View... (Tucker & Snavely 2020)

Typical deep learning pipelines - images go into a 3D CNN, big RGBA 3D volume comes out



Neural Volumes

(Lombardi et al. 2019)

Direct gradient descent to optimize an RGBA volume, regularized by a 3D CNN

