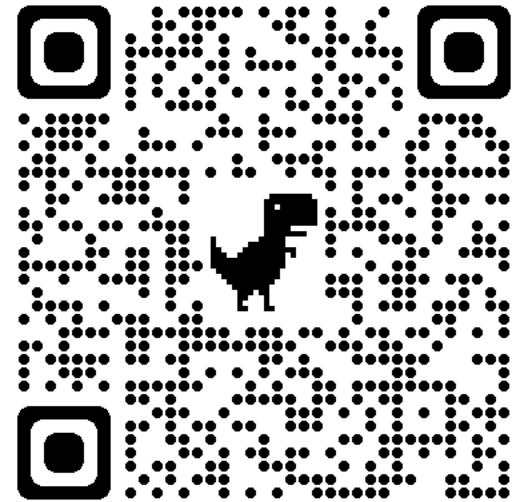# Lecture 22:
# Neural Radiance Fields (NeRFs)

COMP 590/776: Computer Vision
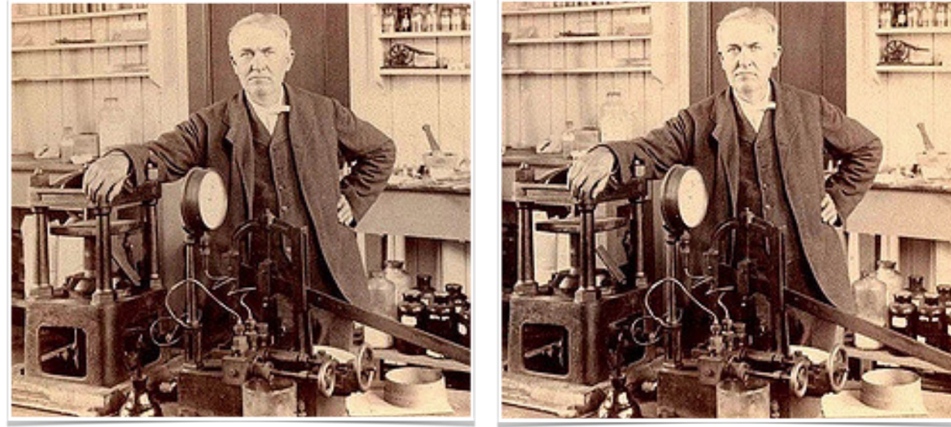
Instructor: Soumyadip (Roni) Sengupta

TA: Mykhailo (Misha) Shvets
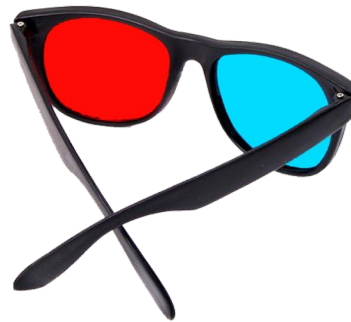
Course Website:
Scan Me!

# Stereo Photography

## Viewing Devices

NeRF (Neural Radiance Field) has revolutionized Computer Vision & Graphics in past 3 years!

Let's look at some of the stunning results it produced!

Given a set of sparse views of an object with known camera poses

Optimize a NeRF model

3D reconstruction viewable from any angle

**NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis**,
Ben Mildenhall*, Pratul Srinivasan*, Matthew Tancik*, Jonathan Barron, Ravi Ramamoorthi, Ren Ng, ECCV 2020.

**NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis**,
Ben Mildenhall*, *Pratul Srinivasan*, Matthew Tancik*, Jonathan Barron, Ravi Ramamoorthi, Ren Ng, ECCV 2020.

10

Block-NeRF: Scalable Large Scene Neural View Synthesis, CVPR 2022.
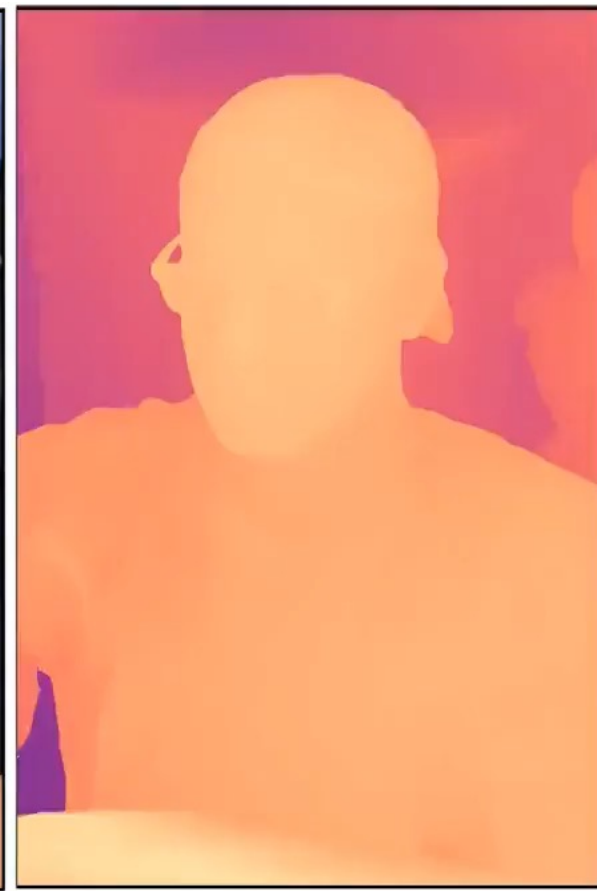
(a) Capture Process  (b) Input  (c) Nerfie  (d) Nerfie Depth

**NeRFies: Deformable Neural Radiance Fields**, Keunhong Park et al., ICCV 2021.

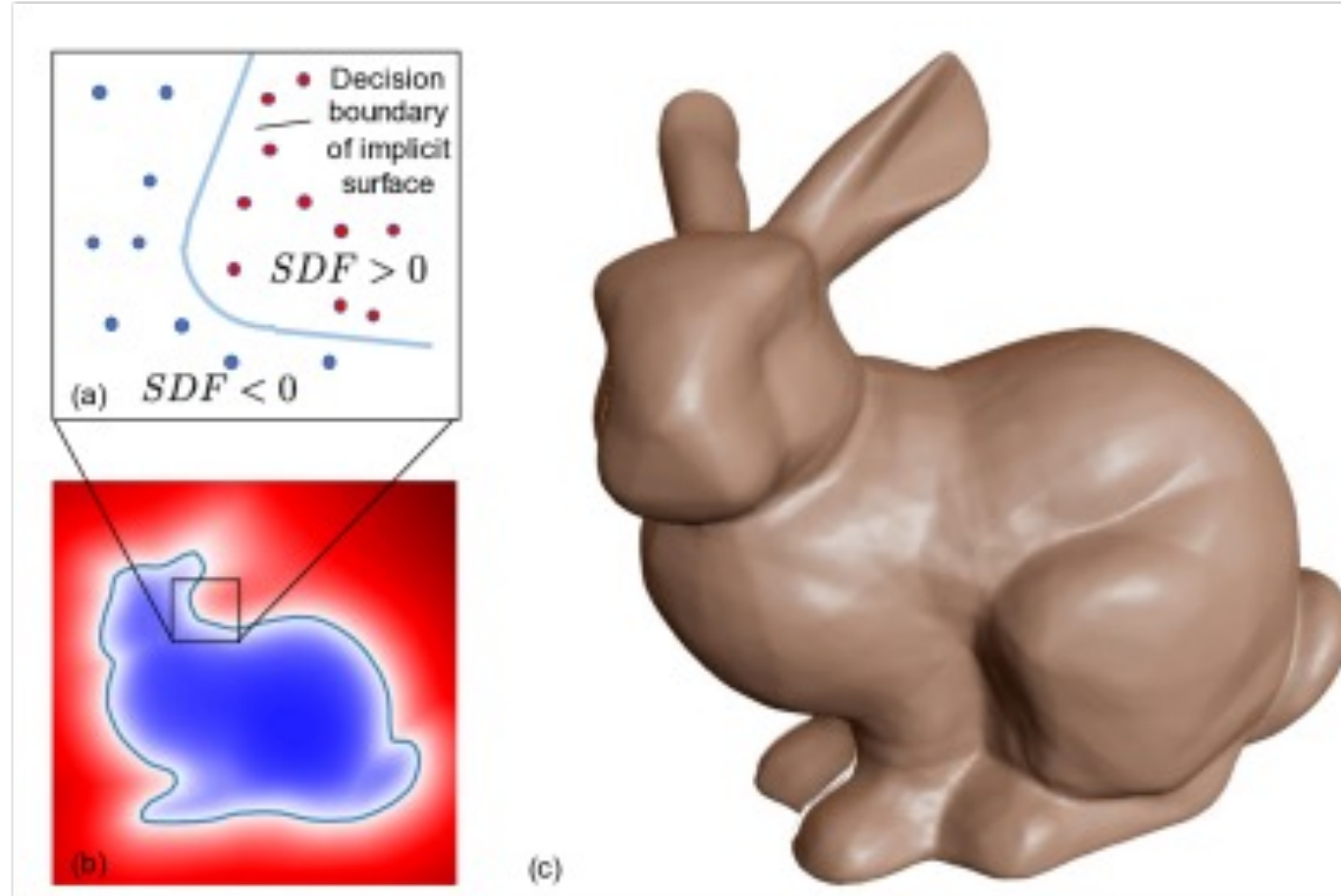Neural 3D Video Synthesis
from Multi-view Video,
Li et al., CVPR 2022

# Surface Representation:
# Signed Distance Function (SDF)
# - implicit representation via level set

SDF(X) = 0, when X is on the surface.
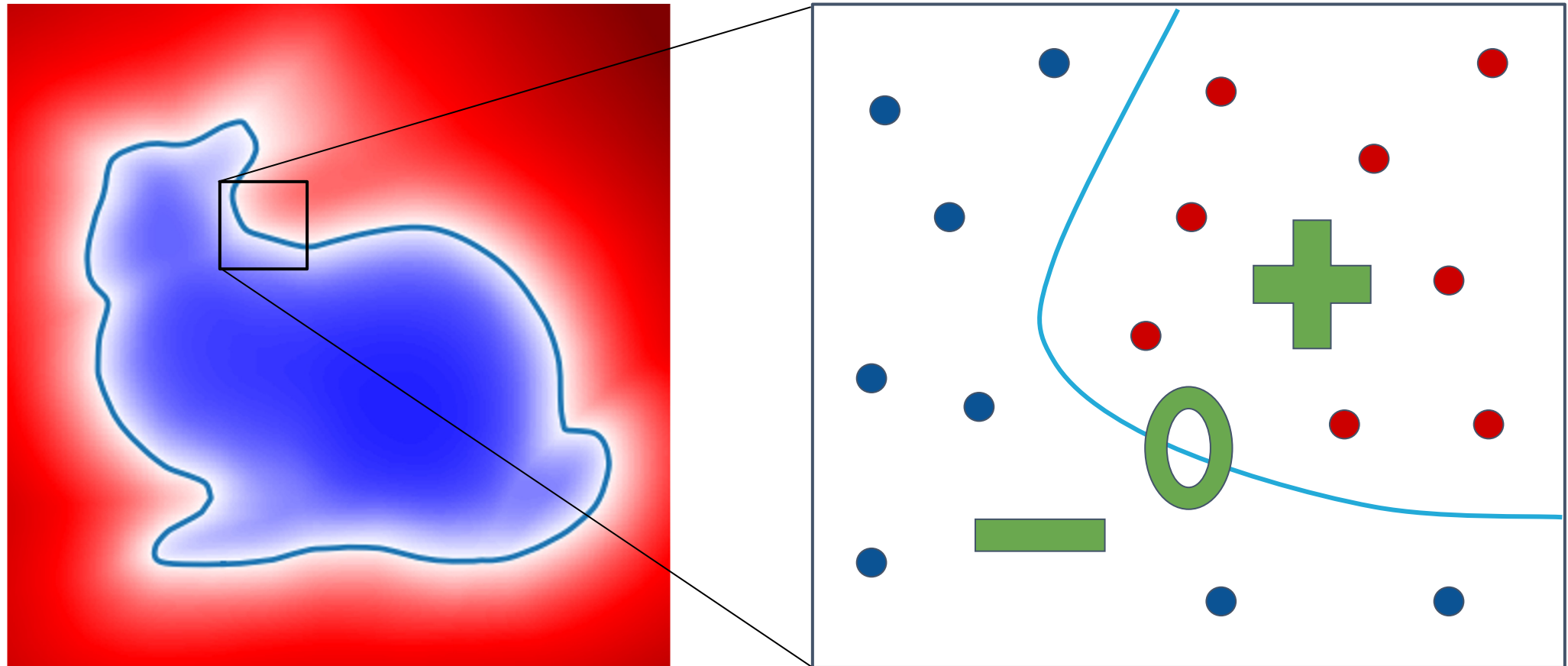
SDF(X) > 0, when X is outside the surface

SDF(X) < 0, when X is inside the surface

Note: SDF is an implicit representation!
Suitable for neural networks but hard to
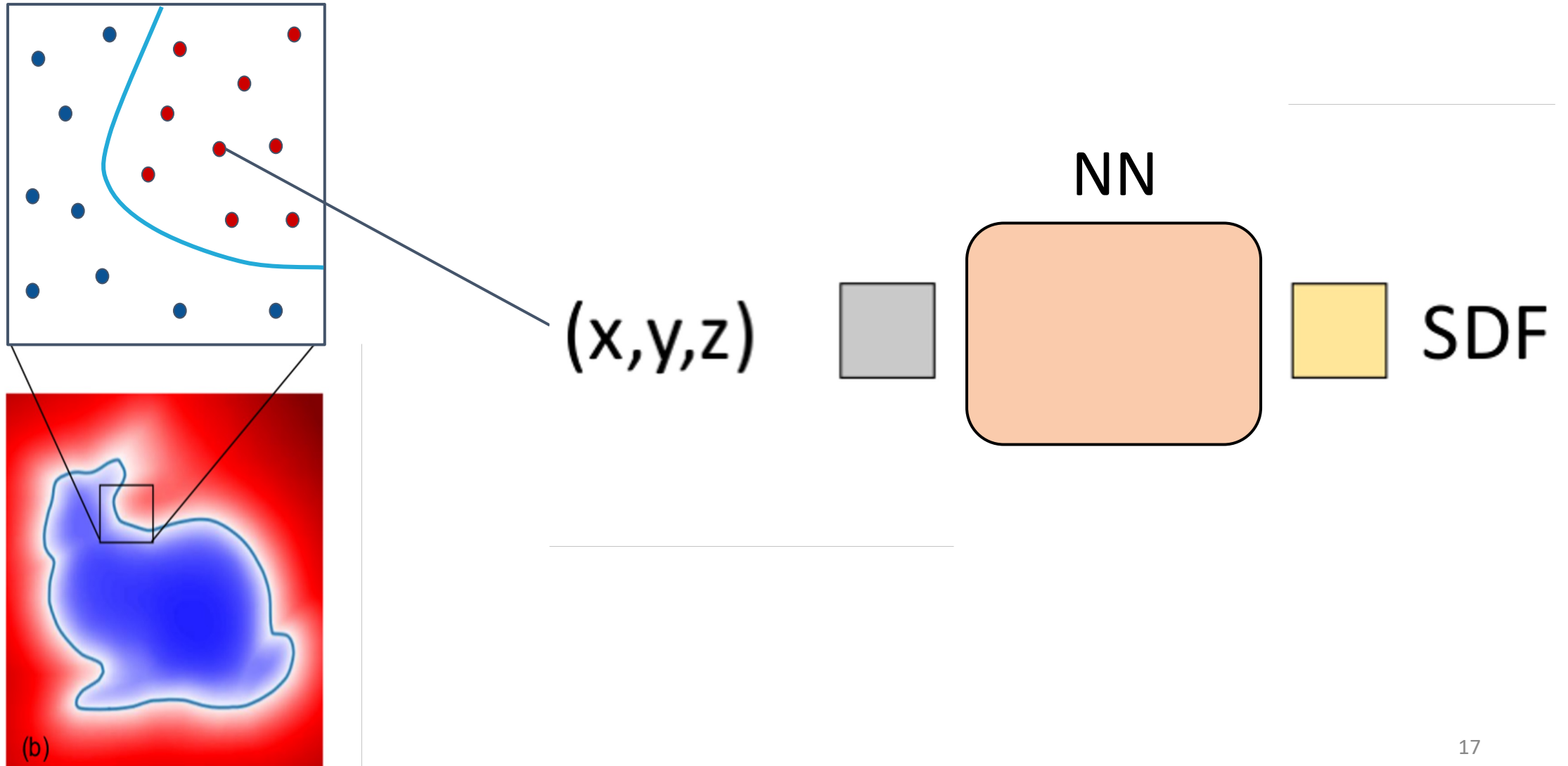import inside existing graphics software.



**Deep SDF: Use a neural network (co-ordinate based MLP) to represent the SDF function.**

# Signed Distance Function

# Regression of Continuous SDF



$(x,y,z)$

NN

SDF

# What is Volume Rendering?

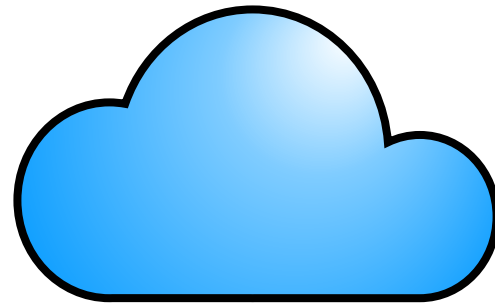- Assume a cloud of tiny colored particles in 3D. Each particle has a RGB color and a density.

- Take a pixel on image plane, and shoot a ray from the camera center, through the pixel and into the 'cloud of tiny colored particles'

- What should be the color for that pixel?

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

Camera

# Volumetric formulation for NeRF



Scene is a cloud of colored fog

# Volumetric formulation for NeRF

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

$t$

Camera

Consider a ray traveling through the scene, and a point at distance $t$ along this ray. We look up its color $\mathbf{c}(t)$, and its opacity (alpha value) $\alpha(t)$

# Volumetric formulation for NeRF



$P[\text{no hits before } t] = T(t)$

$t$

But $t$ may also be blocked by earlier points along the ray. $T(t)$: probability that the ray didn't hit any particles earlier.
$T(t)$ is called "transmittance"

# Volume rendering estimation: integrating color along a ray

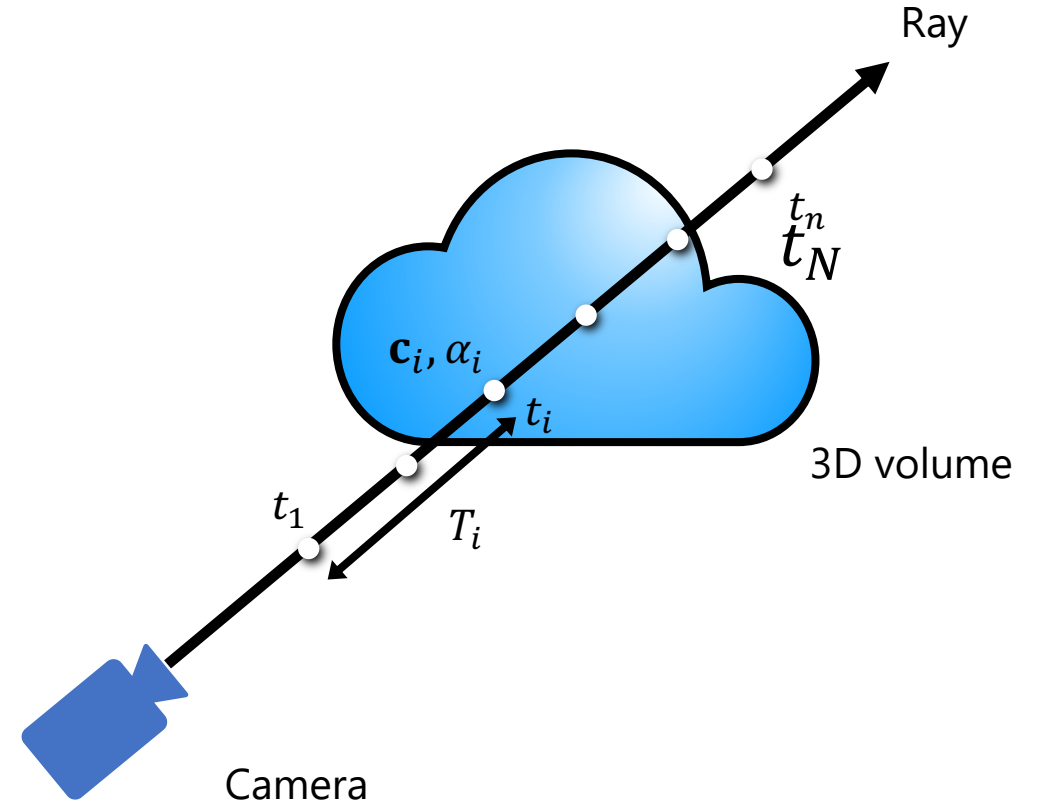Rendering model for ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$:

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

final rendered color along ray

weights

colors

How much light is blocked earlier along ray:

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

Computing the color for a set of rays through the pixels of an image yields a rendered image



Ray

$t_n$
$t_N$

$\mathbf{c}_i, \alpha_i$

$t_i$

3D volume

$t_1$

$T_i$

Camera

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

Slight modification: $\alpha$ is not directly stored in the volume, but instead is derived from a stored volume density sigma ($\sigma$) that is multiplied by the distance between samples delta ($\delta$):

# Volume rendering estimation: integrating color along a ray

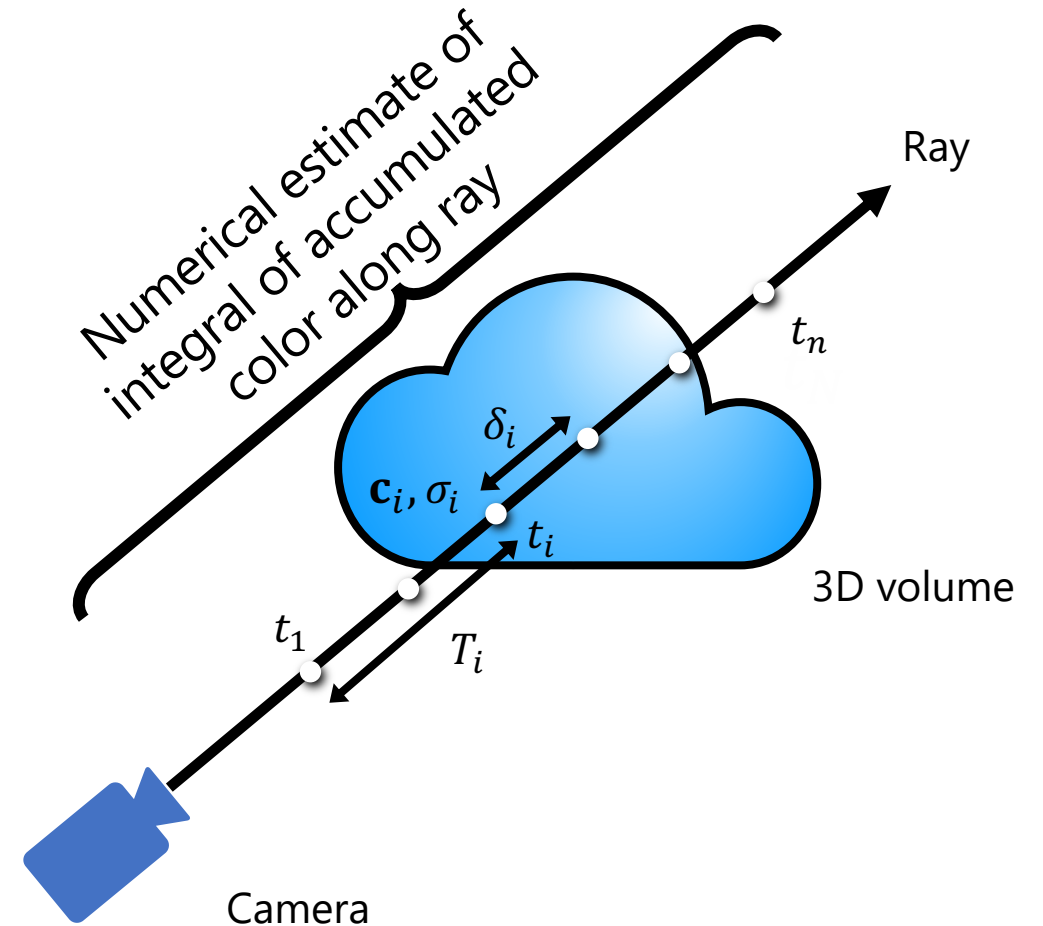Rendering model for ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$:

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

final rendered
color along ray

weights

colors

How much light is blocked earlier along ray:

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

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

Numerical estimate of integral of accumulated color along ray

Ray

$t_n$

$\delta_i$

$\mathbf{c}_i, \sigma_i$

$t_i$

3D volume

$t_1$

$T_i$

Camera

# Volume rendering estimation: integrating color along a ray

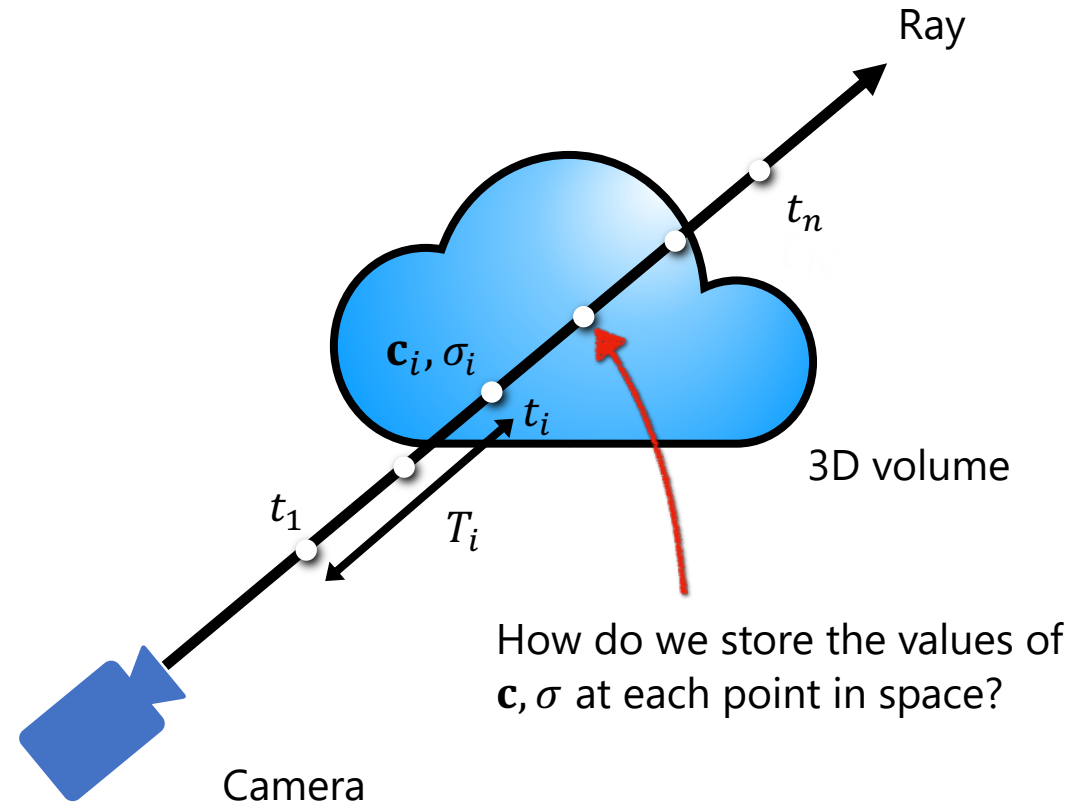Rendering model for ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$:

$$\mathbf{c} \approx \sum_{i=1}^{n} T_i \alpha_i \mathbf{c}_i$$
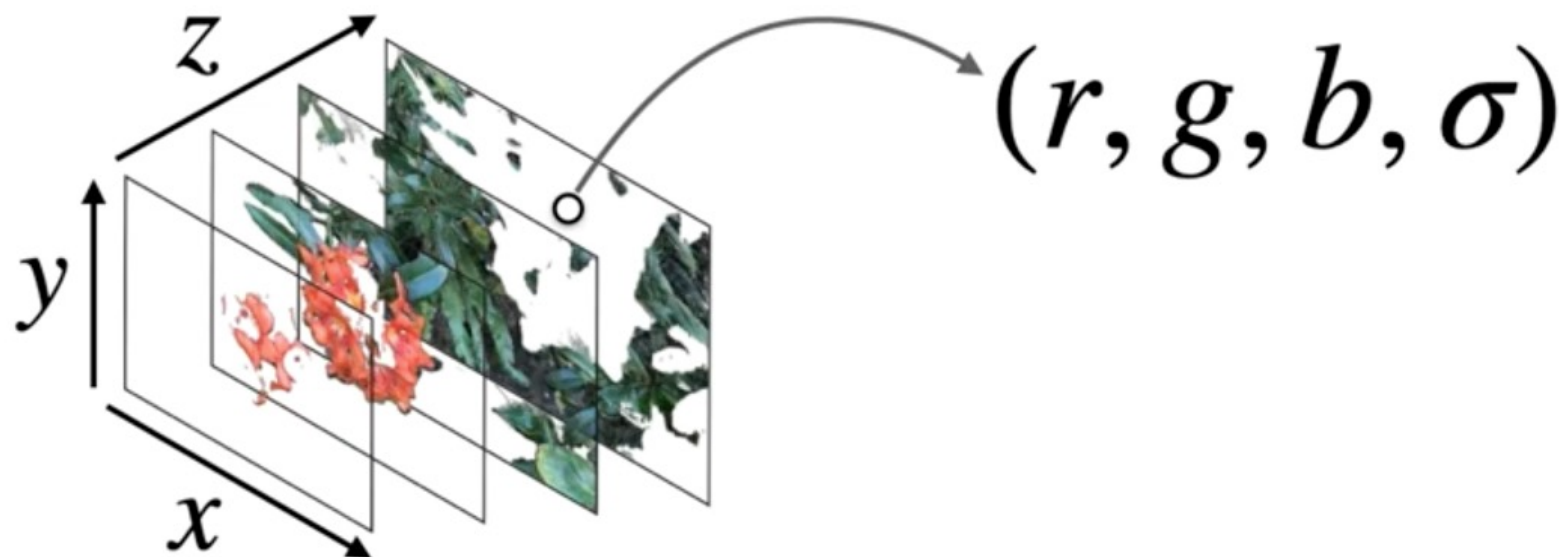
final rendered color along ray

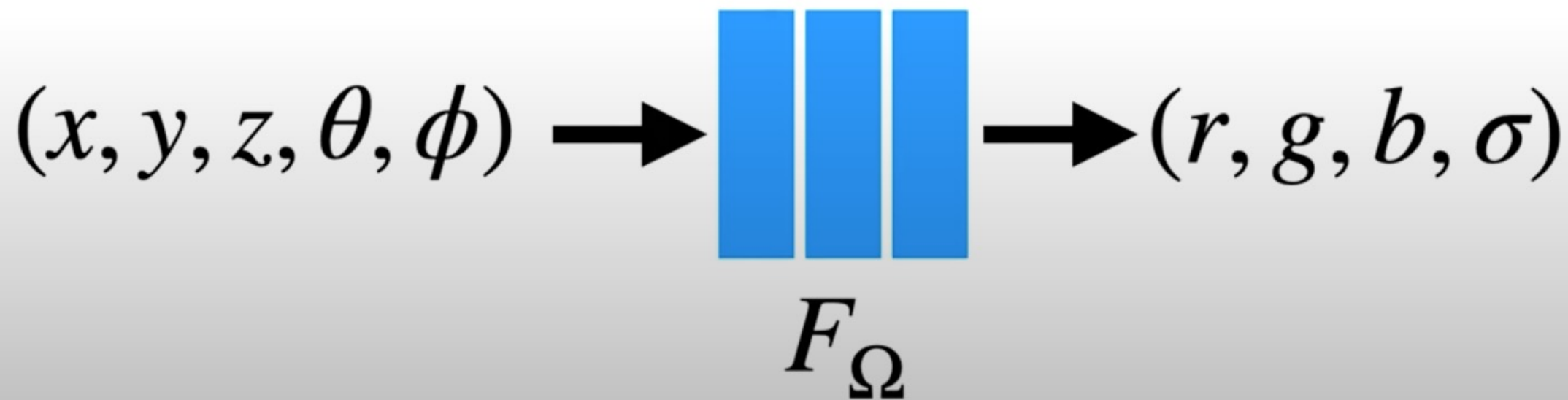weights

colors

How much light is blocked earlier along ray:

$$T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$$
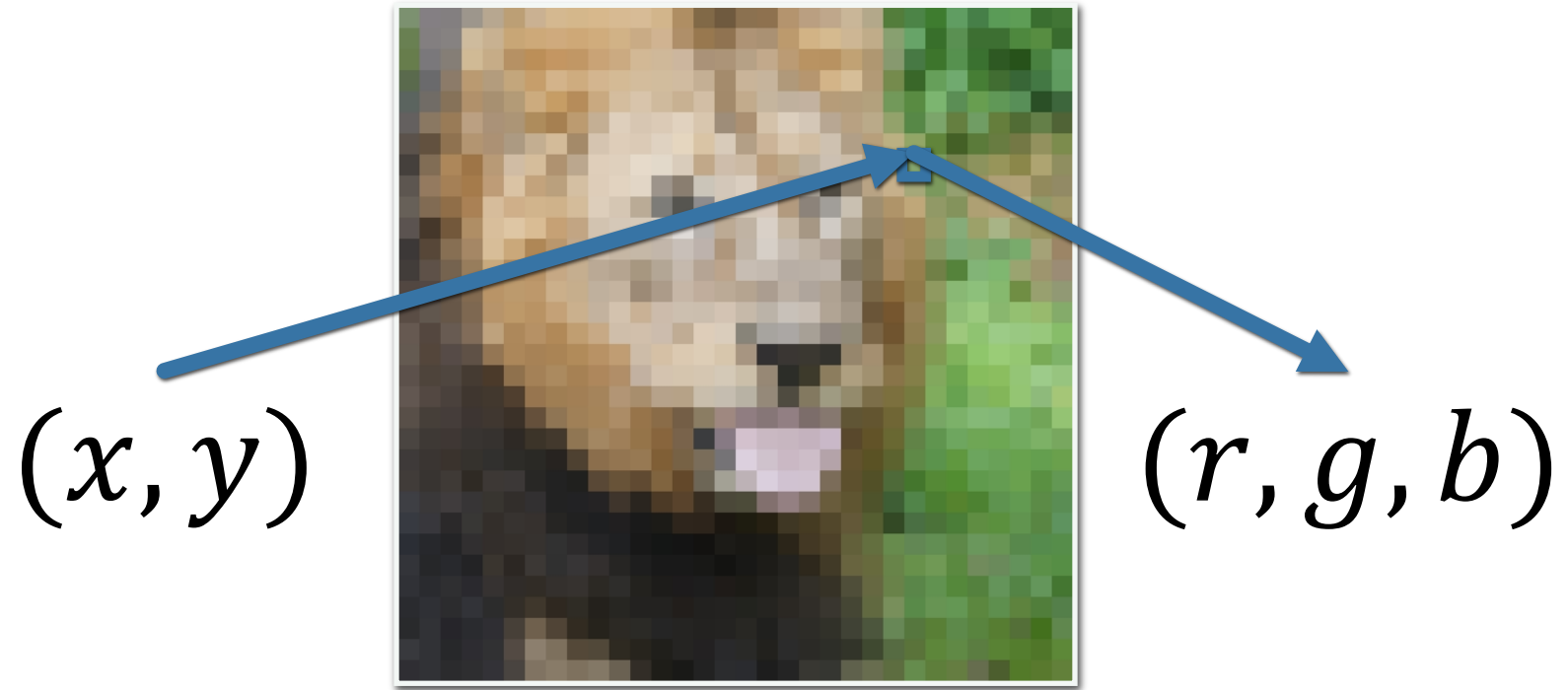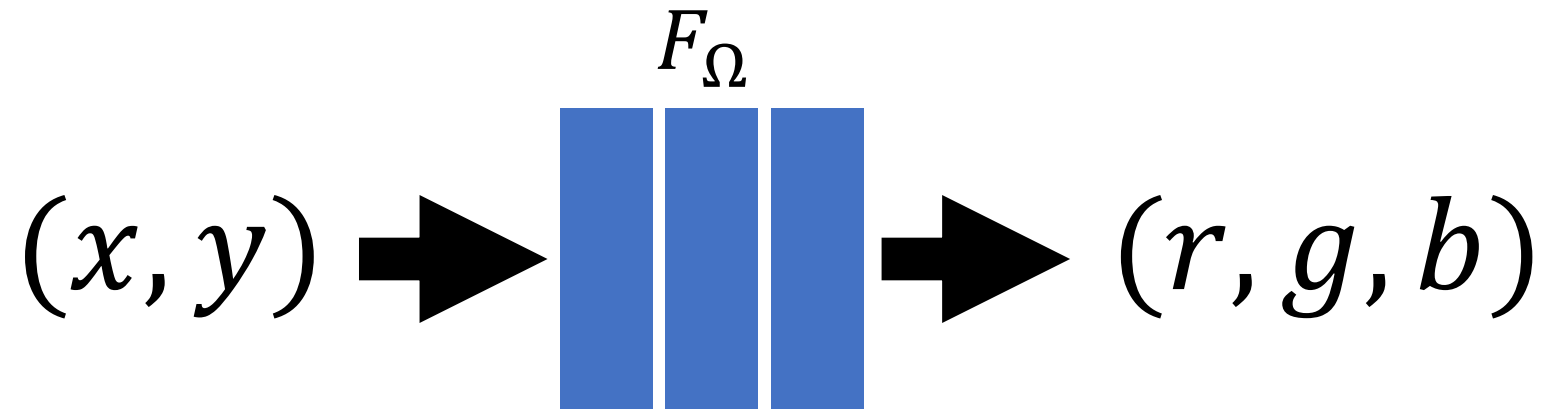
$$\alpha_i = 1 - \exp(-\sigma_i \delta_i)$$

Ray

$t_n$

$\mathbf{c}_i, \sigma_i$

$t_i$

3D volume

$t_1$

$T_i$

How do we store the values of $\mathbf{c}, \sigma$ at each point in space?

Camera

$(r, g, b, \sigma)$

versus

$(x, y, z, \theta, \phi) \longrightarrow F_\Omega \longrightarrow (r, g, b, \sigma)$

Toy problem: storing 2D image data



$(x, y)$

$(r, g, b)$

Usually we store an image as a
2D grid of RGB color values

# Toy problem: storing 2D image data

$$F_{\Omega}$$

$$(x, y) \Rightarrow \quad \blacksquare\blacksquare\blacksquare \quad \Rightarrow (r, g, b)$$

What if we train a simple fully-connected network (MLP) to do this instead?

# Naive approach fails!


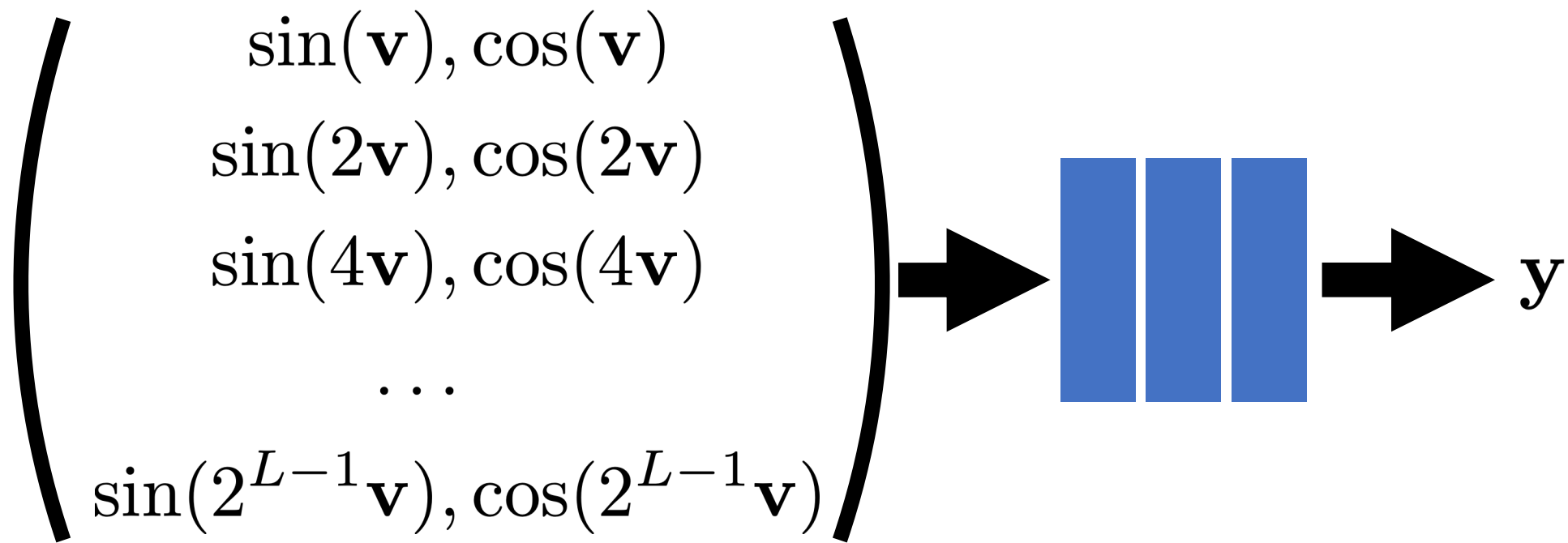Ground truth image


Neural network output fit
with gradient descent

28

# Problem:

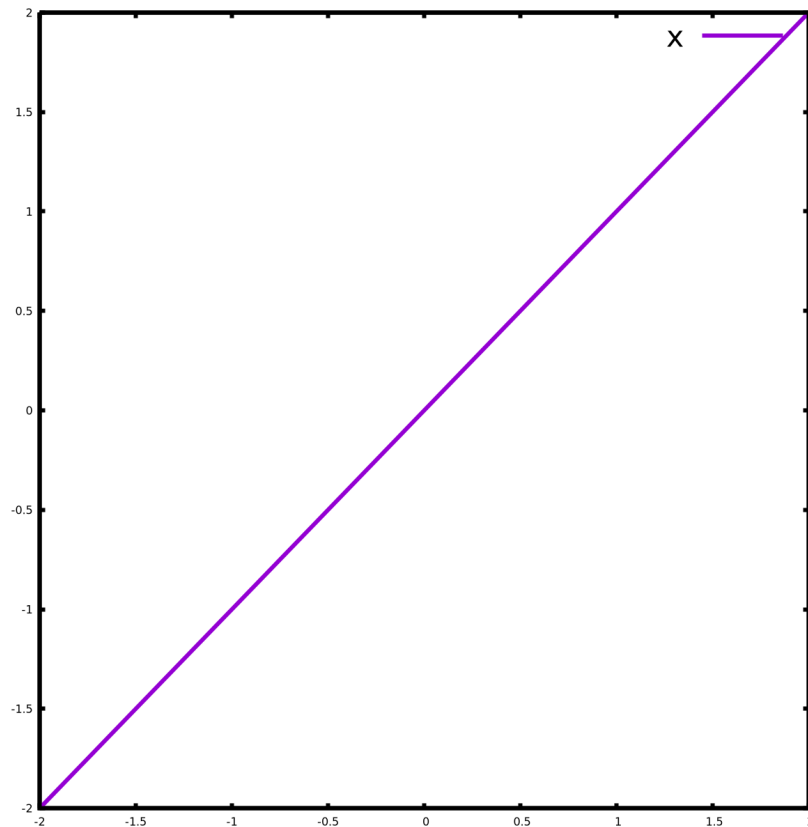- "Standard" coordinate-based MLPs cannot represent high frequency functions.

# Solution:

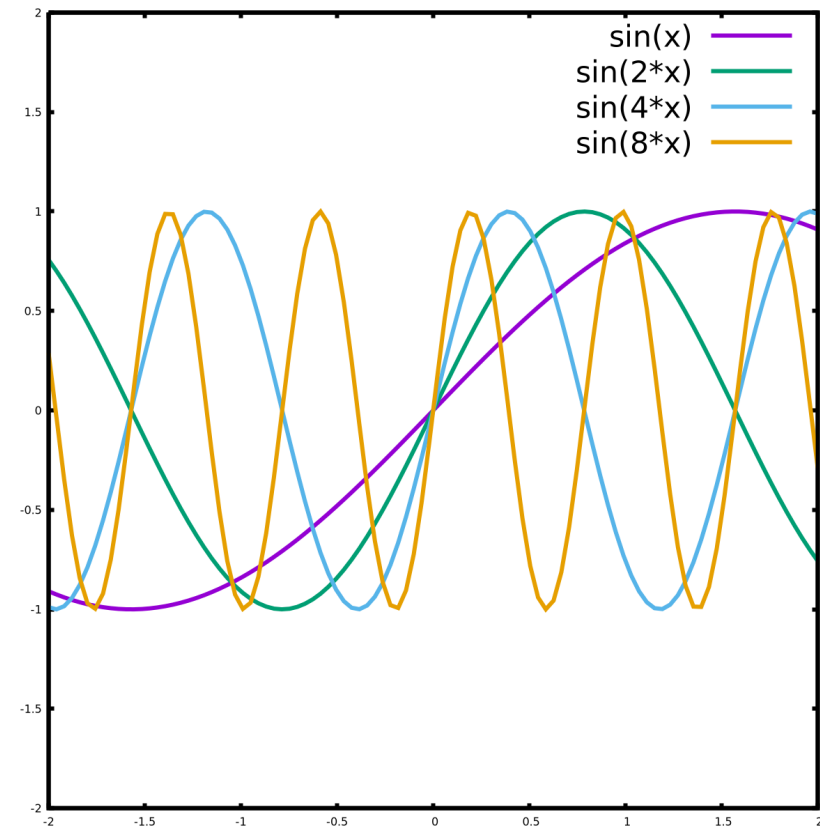- Pass input coordinates through a high frequency mapping first.

Example mapping: "positional encoding"

$$\mathbf{v} \rightarrow \boxed{\phantom{|||}} \rightarrow \mathbf{y}$$

$$\begin{pmatrix} \sin(\mathbf{v}), \cos(\mathbf{v}) \\ \sin(2\mathbf{v}), \cos(2\mathbf{v}) \\ \sin(4\mathbf{v}), \cos(4\mathbf{v}) \\ \dots \\ \sin(2^{L-1}\mathbf{v}), \cos(2^{L-1}\mathbf{v}) \end{pmatrix} \rightarrow \boxed{\phantom{|||}} \rightarrow \mathbf{y}$$

# Positional encoding



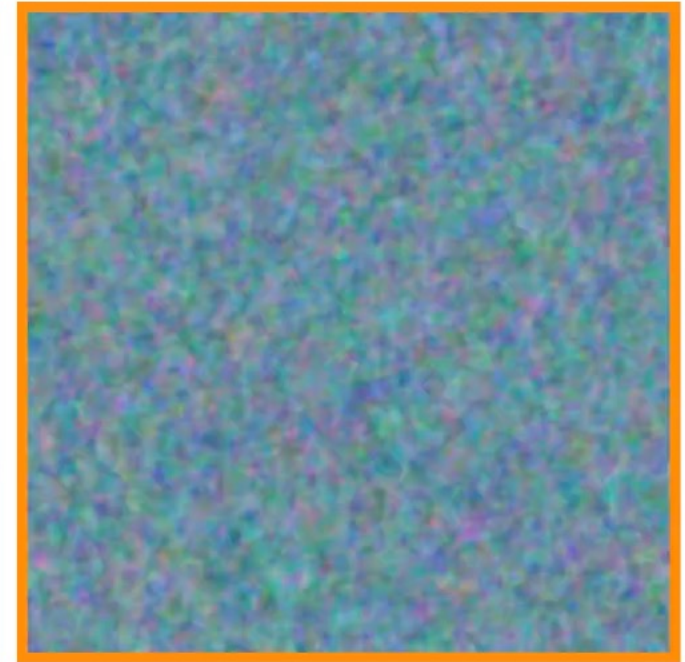Raw encoding of a number **x**

"Positional encoding" of a number **x**
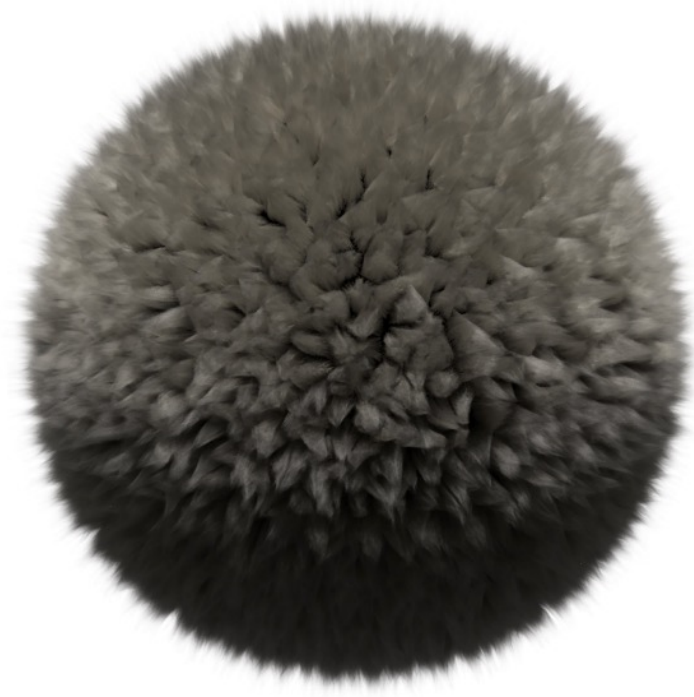
# Problem solved!



Ground truth image

Neural network output without high frequency mapping

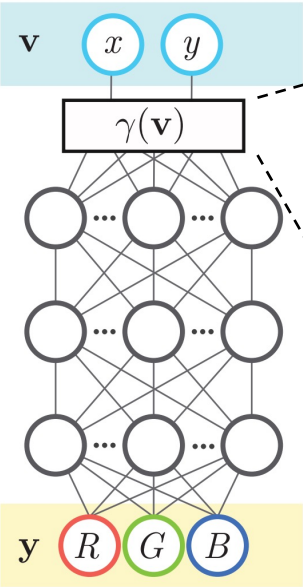Neural network output with high frequency mapping

32

32

# Network Architecture: Overcoming Spectral Bias



[Baatz et al. 2021]

**The signals we want are high frequency!**

# Network Architecture: Input Encodings

**v** $x$ $y$

$\gamma(\mathbf{v})$

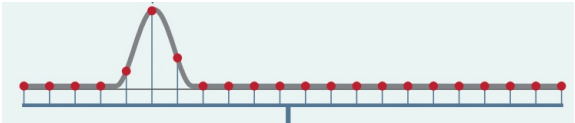**y** $R$ $G$ $B$

## Random Fourier Encodings

[Tancik et al. 2020]

$$\gamma(\mathbf{v}) = [\cos(2\pi\mathbf{B}\mathbf{v}), \sin(2\pi\mathbf{B}\mathbf{v})]^{\mathrm{T}}$$

Non-axis aligned sine embeddings
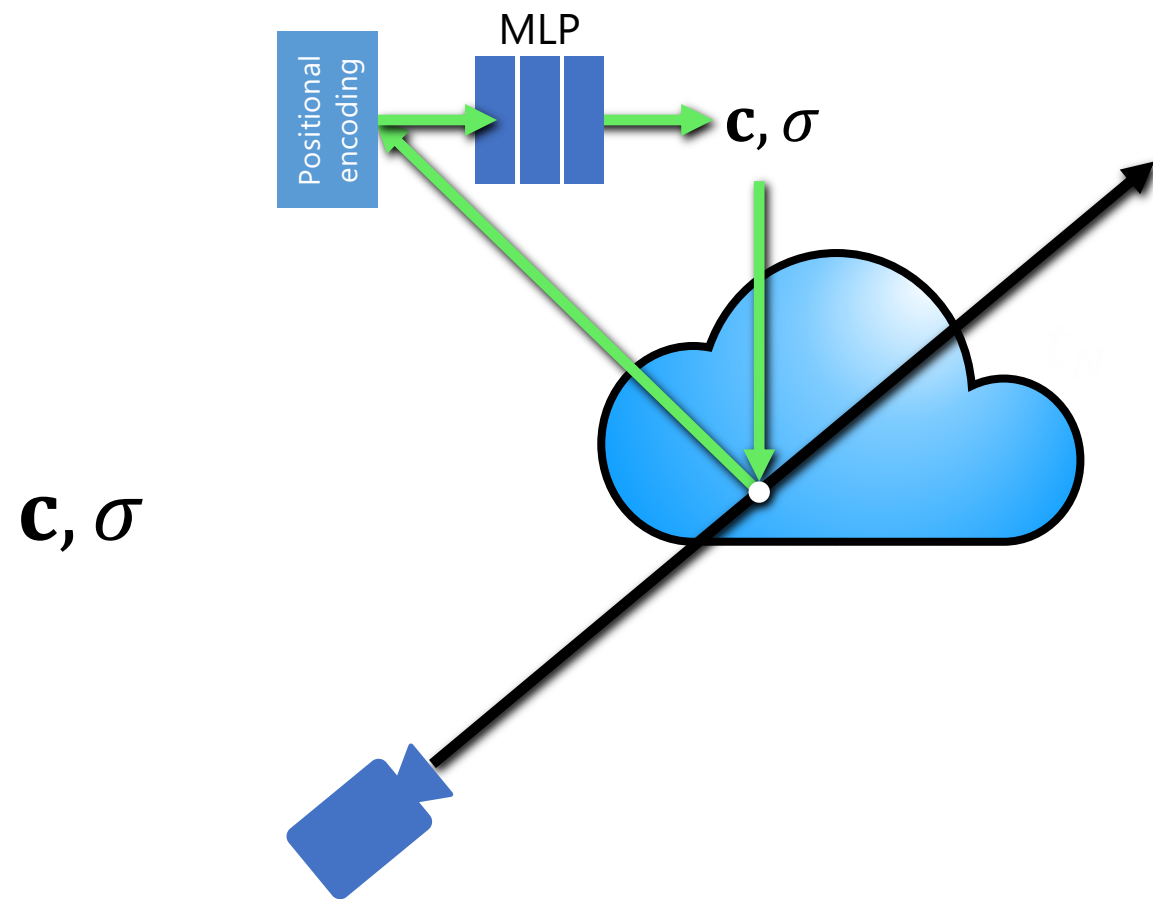
## One-blob Encodings

[Müller et al. 2020]
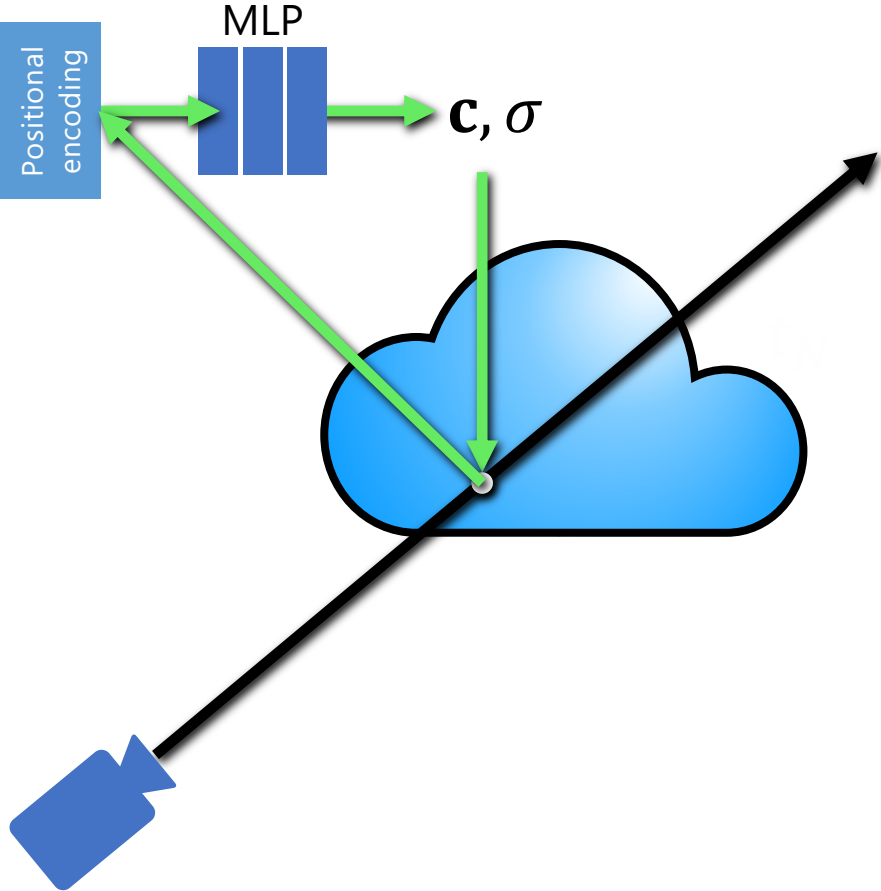
## Super Gaussian Encodings

[Ramasinghe et al. 2021]

$$\Phi(\boldsymbol{x}) = [\phi_1(\boldsymbol{x}), \phi_2(\boldsymbol{x}), \ldots, \phi_D(\boldsymbol{x})]^T,$$

$$\left[ e^{-\frac{(\boldsymbol{x}\cdot\alpha - t_i)^2}{2\sigma_{\boldsymbol{x}}^2}} \right]^b$$

Gaussian embeddings

# NeRF = volume rendering + coordinate-based network

MLP

Positional encoding

$\mathbf{c}, \sigma$

$\mathbf{c}, \sigma$

**c**, $\sigma$

Positional encoding

MLP

**c**, $\sigma$

$\mathbf{c}, \sigma$

Positional encoding

MLP

$\mathbf{c}_0, \sigma_0$

Positional encoding

MLP

$\mathbf{c}_1, \sigma_1$

$\mathbf{c}, \sigma$

Positional encoding

MLP

$\mathbf{c}_3, \sigma_3$

$\mathbf{c}, \sigma$

Positional encoding

$\mathbf{c}_n, \sigma_n$
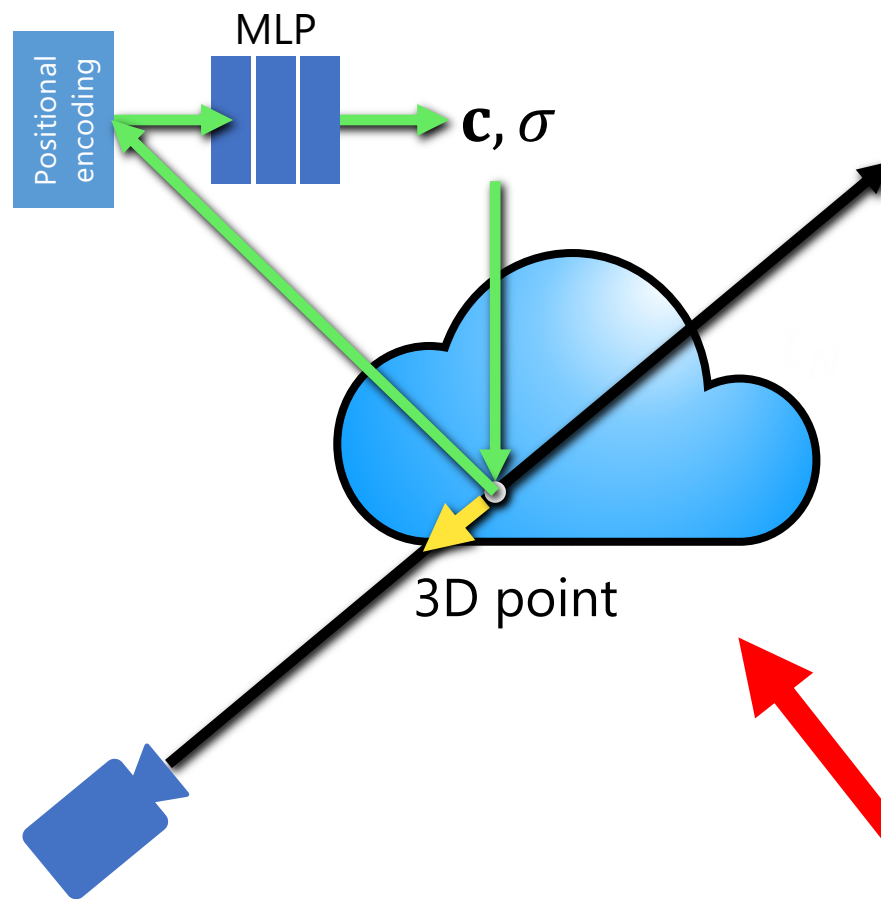
MLP

MLP

Positional encoding

$\mathbf{c}, \sigma$

3D point
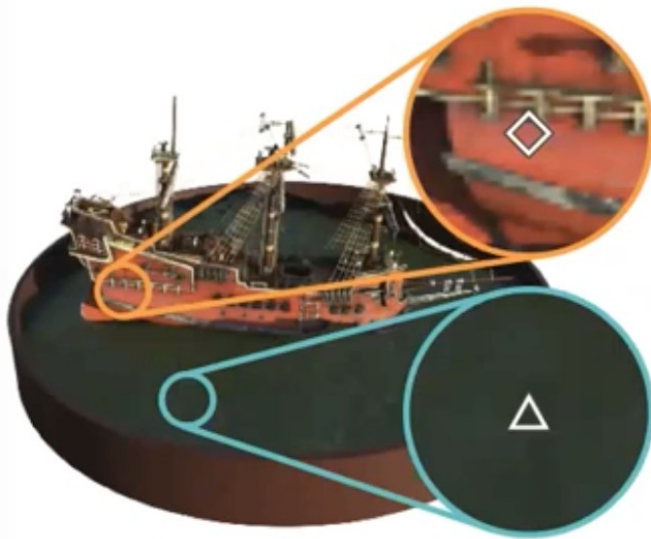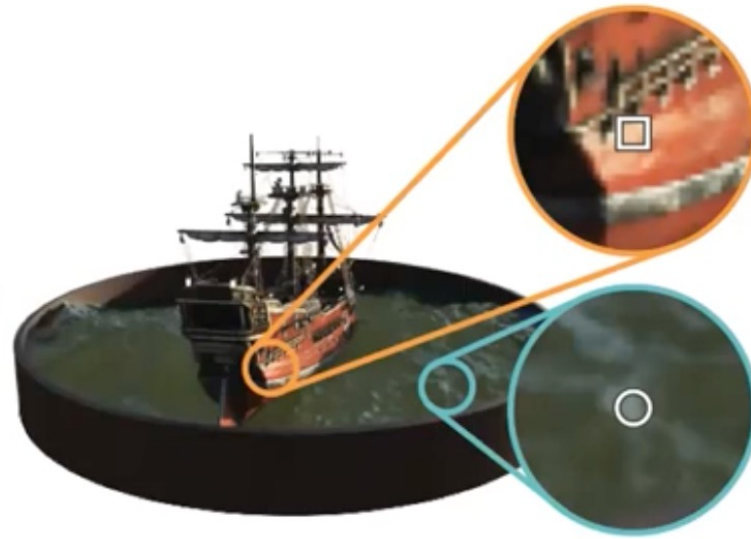
Include the ray direction in the input to the MLP → allows for capturing and rendering view-dependent effects (e.g., shiny surfaces)

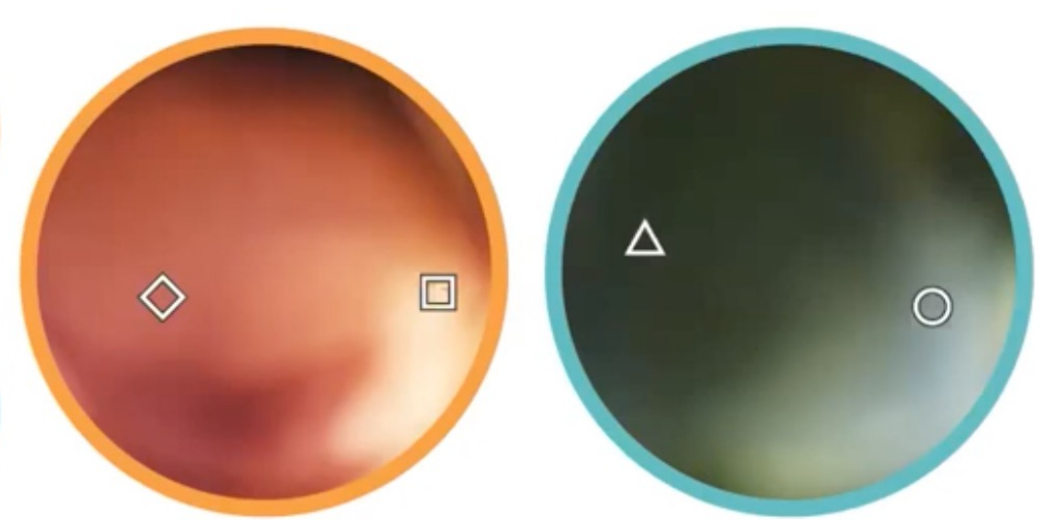# Modeling view dependent effects



(a) View 1        (b) View 2        (c) Radiance Distributions
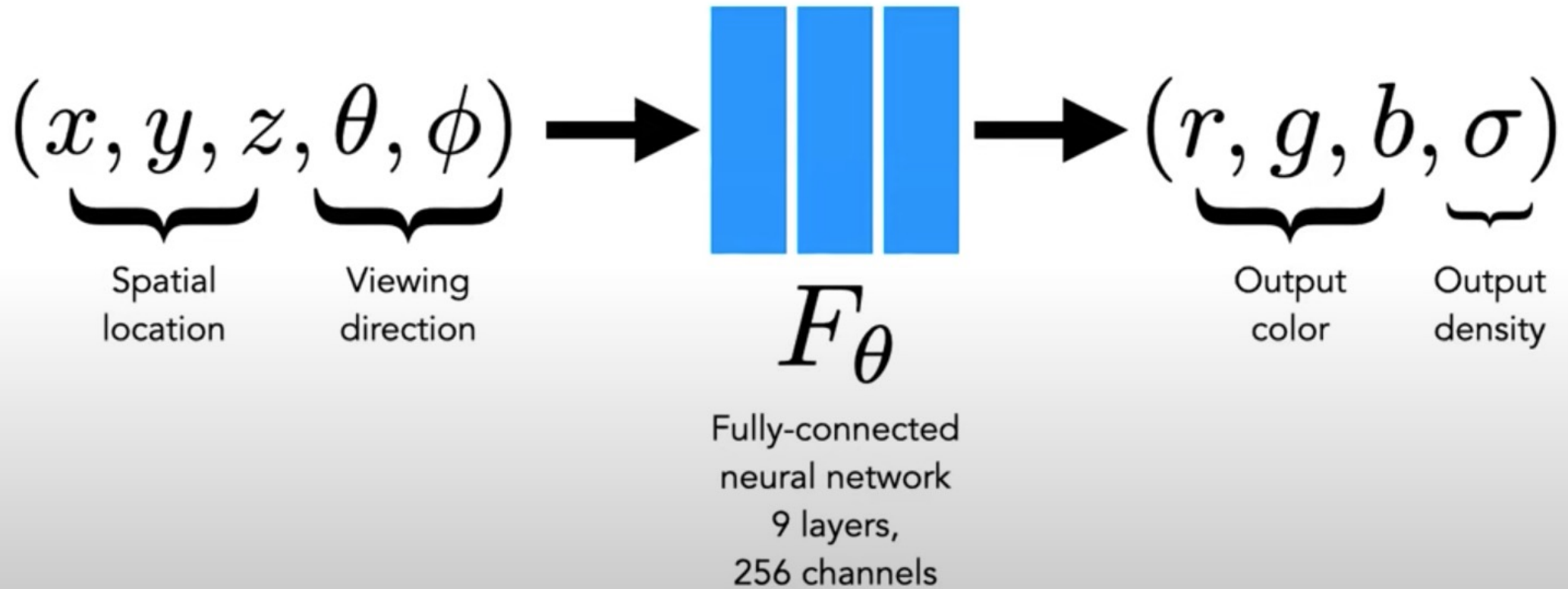
$$(x, y, z, \theta, \phi) \rightarrow \boxed{F_\theta} \rightarrow (r, g, b, \sigma)$$

Spatial location | Viewing direction

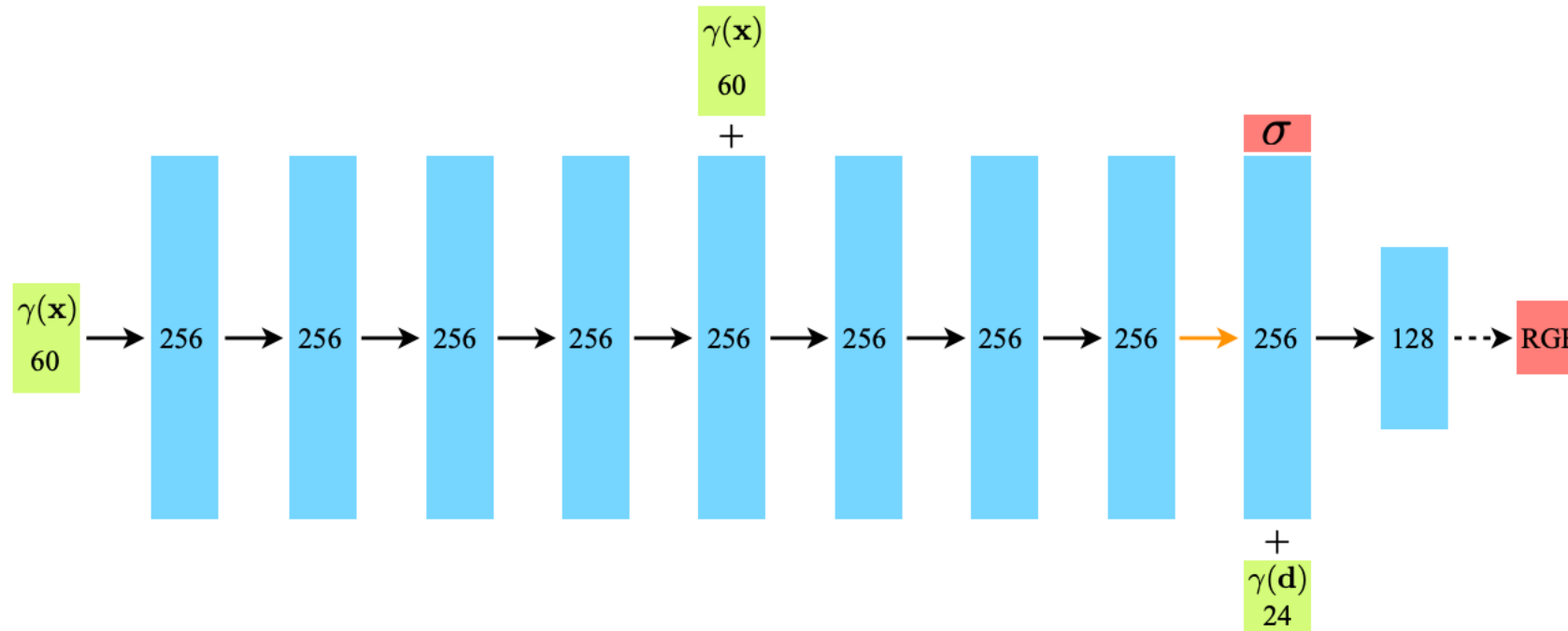$F_\theta$

Fully-connected neural network
9 layers,
256 channels

Output color | Output density

# What do we learn in NeRF?

$$(x, y, z, \theta, \phi) \rightarrow \boxed{\,\,\, \text{III} \,\,\,} \rightarrow (r, g, b, \sigma)$$

Spatial location     Viewing direction

$F_\theta$

Fully-connected neural network
9 layers,
256 channels

Output color     Output density

# DeepSDF Extensions: NeRF

- Coordinate-based modeling of RGB and Densities Instead of SDFs



Mildenhall et al. 2020

# Training NeRFs

Volume rendering of MLP colors/densities

Ground truth image

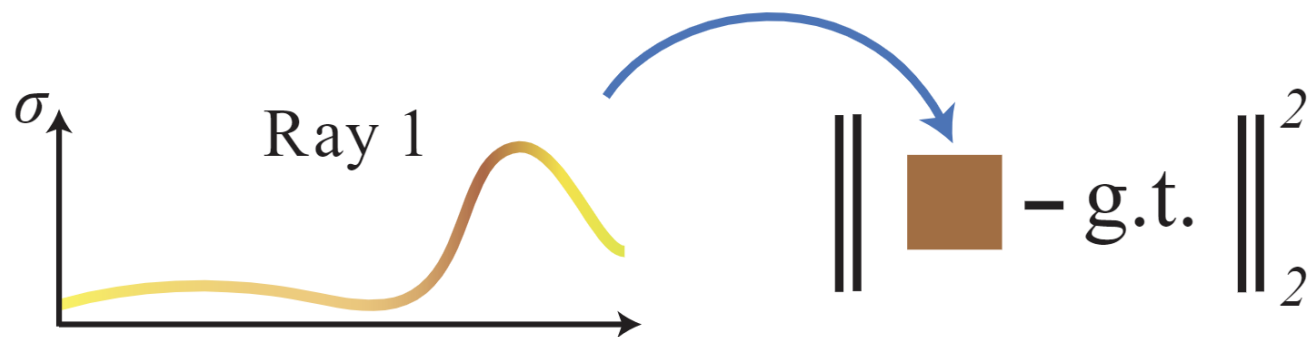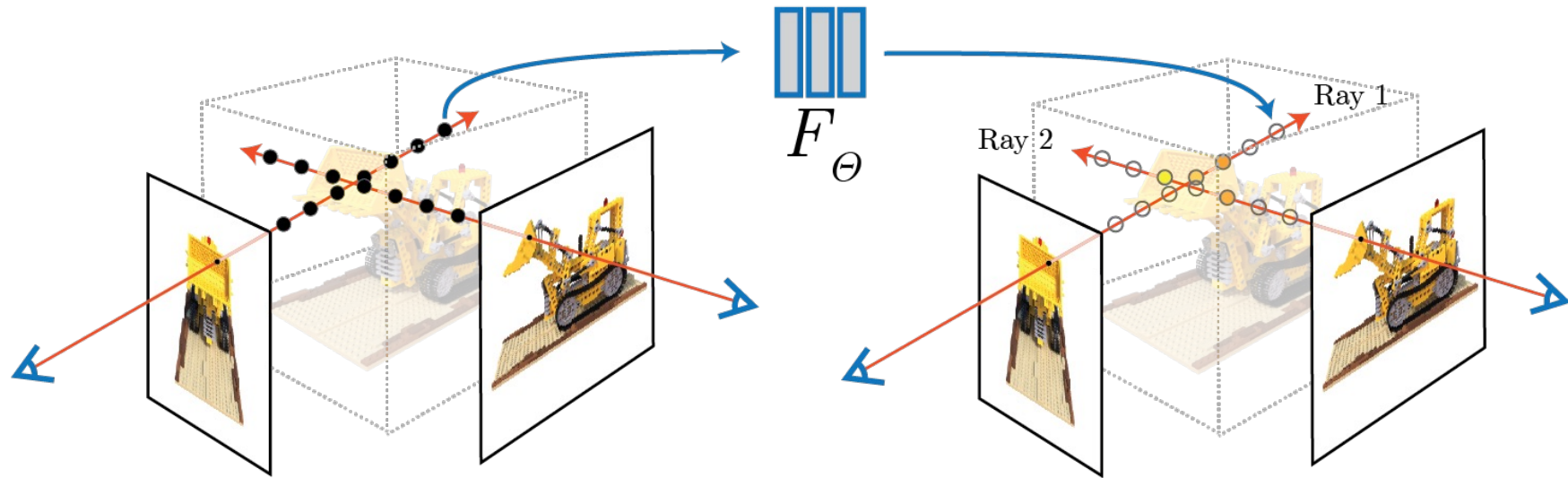$$\nabla \left\| \phantom{x} - \phantom{x} \right\|^2$$

$$\min_\theta \sum_i \| \operatorname{render}_i(F_\theta) - I_i \|^2$$

# Importance Sampling

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

treat weights as probability distribution for new samples

Ray

3D volume

Camera

# NeRF encodes convincing view-dependent effects using directional dependence

# Building 3D models from NeRFs

Apply marching cubes algorithm on NeRF predicted volume density ($\sigma$)
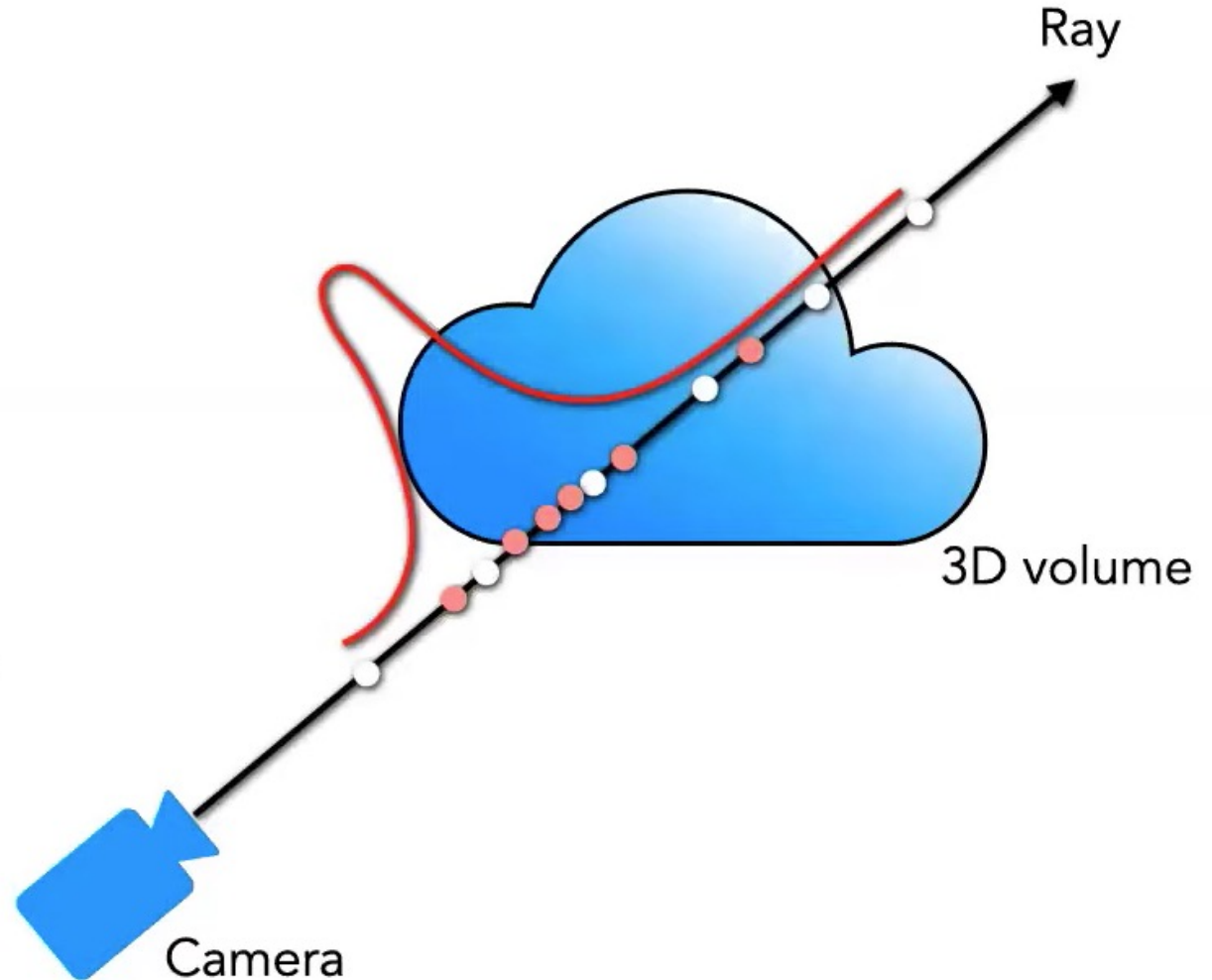
# Summary

- Represent the scene as volumetric colored "fog"
- Store the fog color and density at each point as an MLP mapping 3D position (x, y, z) to color c and density $\sigma$
- Render image by shooting a ray through the fog for each pixel
- Optimize MLP parameters by rendering to a set of known viewpoints and comparing to ground truth images

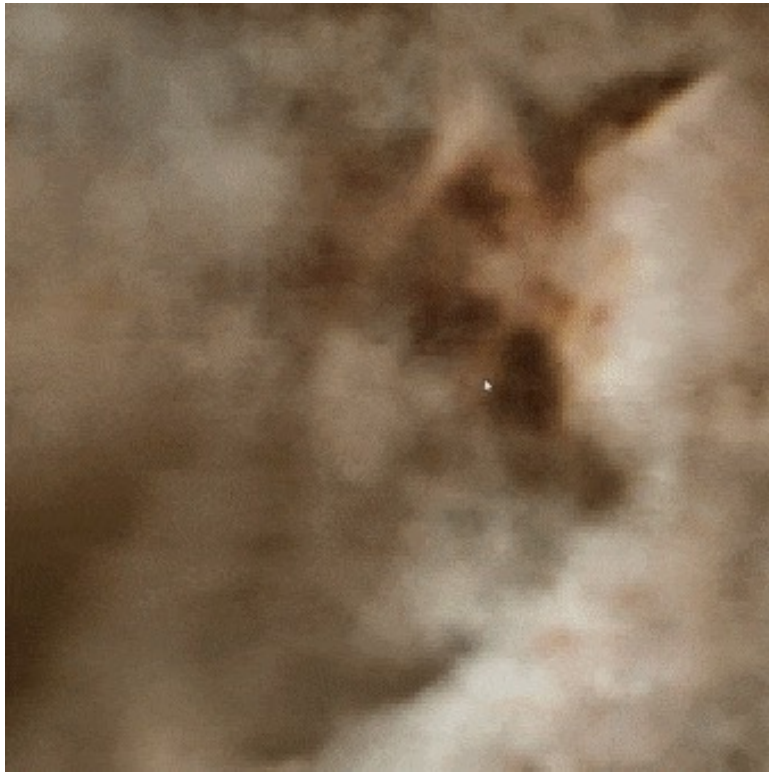# Key limitations of the original NeRF

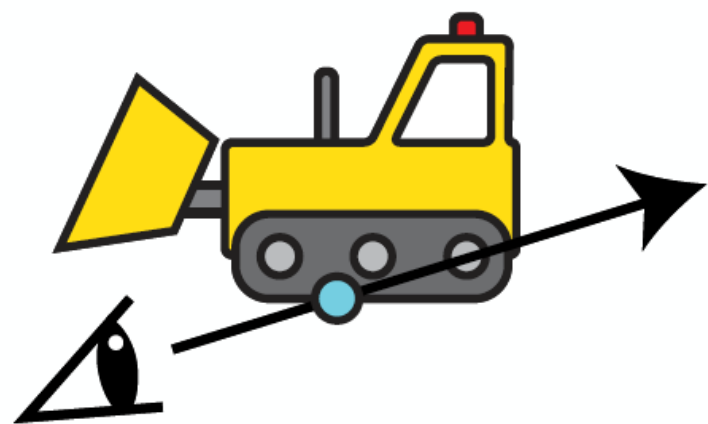- Very slow in training and inference
- Requires Ground-Truth poses
- Do not generalize to new scenes

# Key limitations of the original NeRF
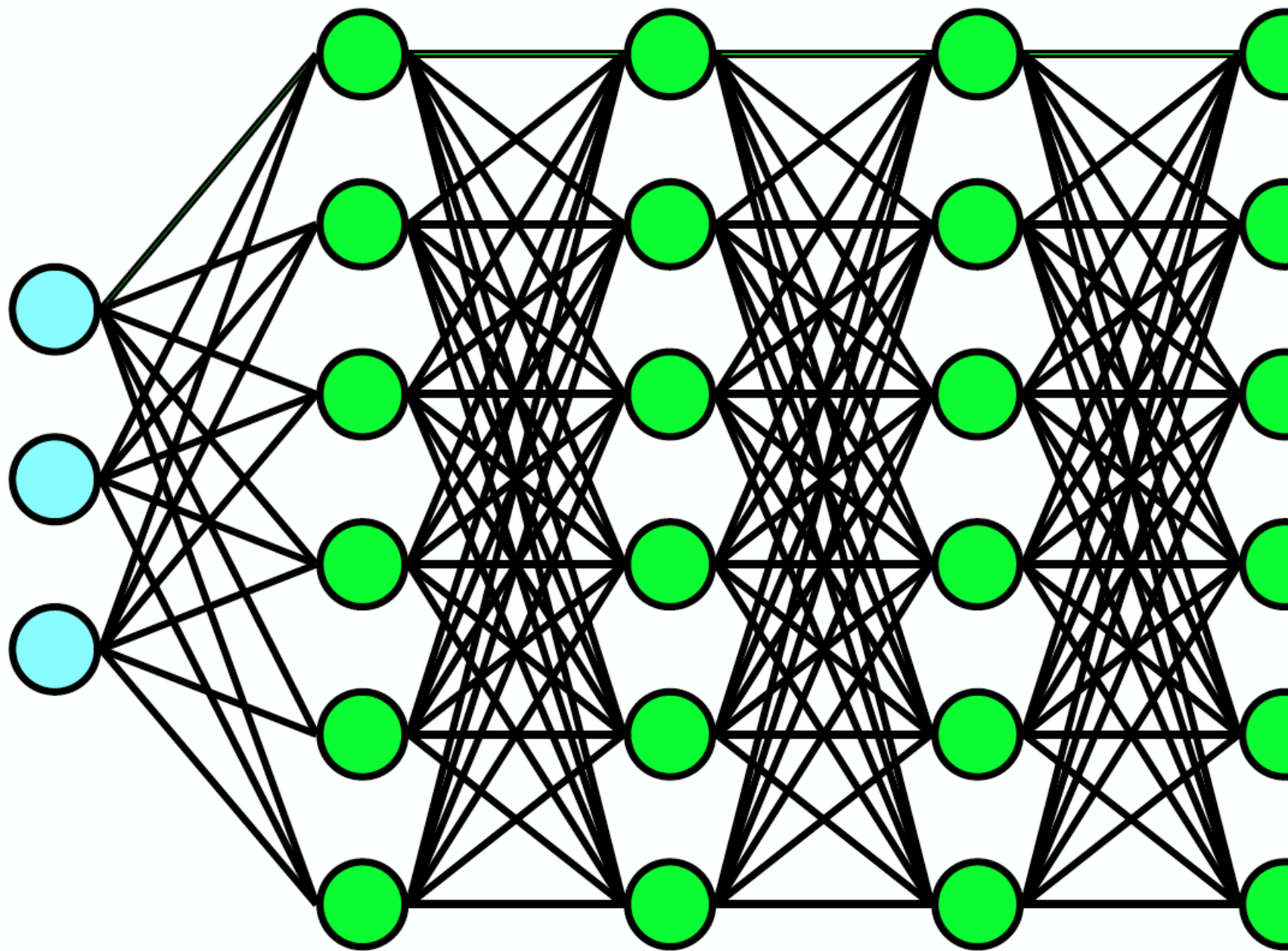
- <span style="color:red">Very slow in training and inference</span>
- Requires Ground-Truth poses
- Do not generalize to new scenes

# Instant NGP: Superfast training and inference with NeRF using multi-resolution hash-table

Ray Query Point

Huge Neural Network 🙁

4

# Hybrid representation



Ray Query Point          Feature Grid Interpolation          Tiny Neural Network 😊
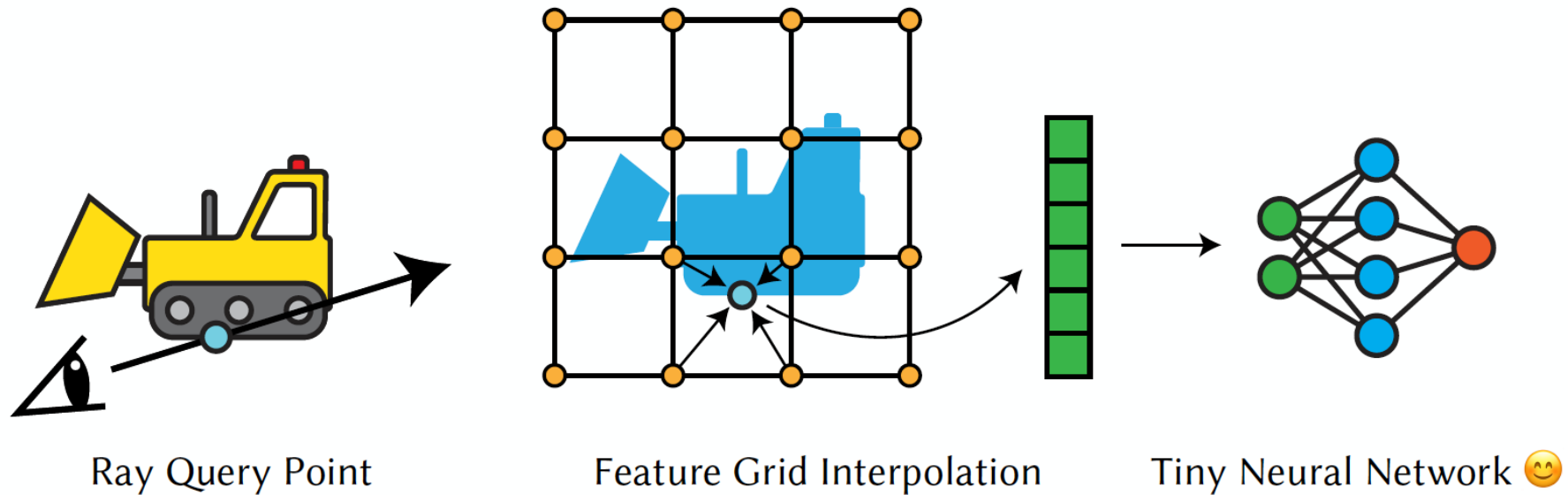
Features:

- are also parameters that can be updated while training the NeRF. (slight increase in memory, significantly faster training & inference)

- are individual NeRFs trained on a small section of a scene (for large city-size scene)

- are priors obtained from ConvNets, e.g. VGG-features (used for generalization)

# Hybrid representation: It's all about Data Structures!



Ray Query Point            Feature Grid Interpolation            Tiny Neural Network 😊

Why hybrid representation?

- Reduce the size of neural network -> fast inference & rendering.
- Helps in rendering large scale scenes.
- Helps in generalization.

# Uniform Grids



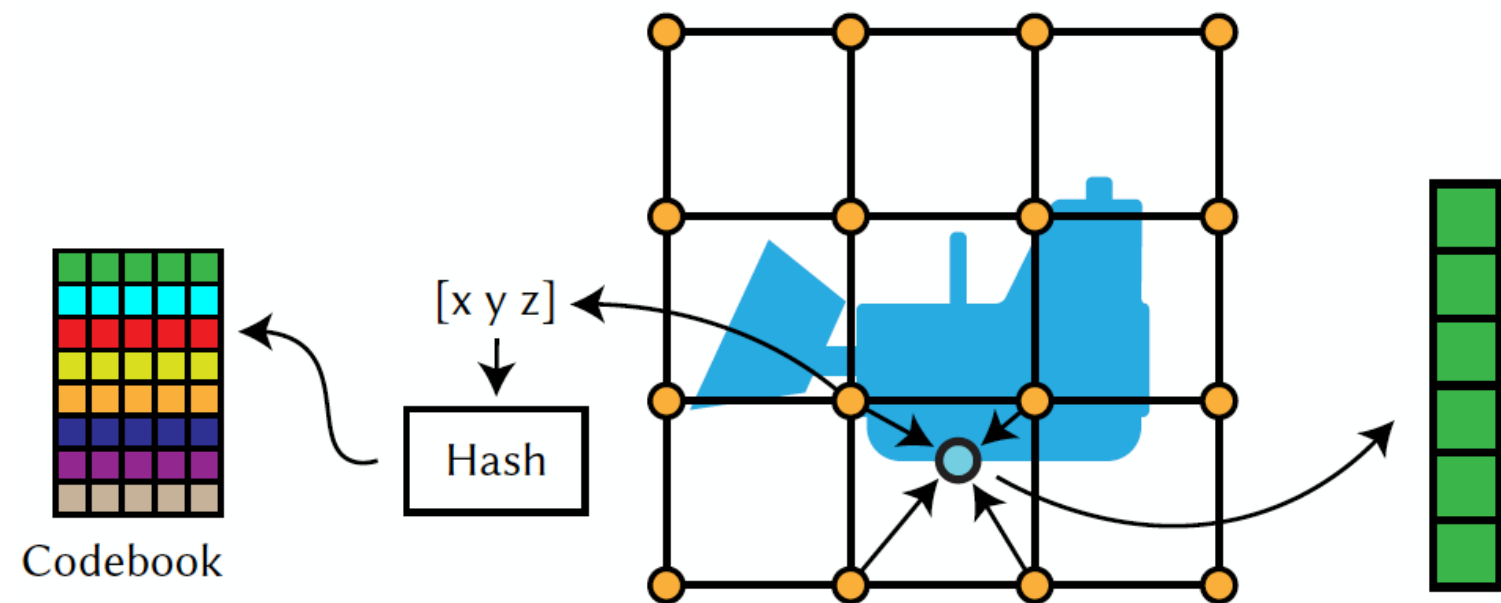[PIFu (Saito et al.), Neural Volumes (Lombardi et al.), etc]

Pros:

- Easy to implement
- Algorithmically fast access [O(1)]
- Established operations like convolutions
- Simple topology

Cons:

- Expensive in memory and bandwidth
- Limited by Nyquist
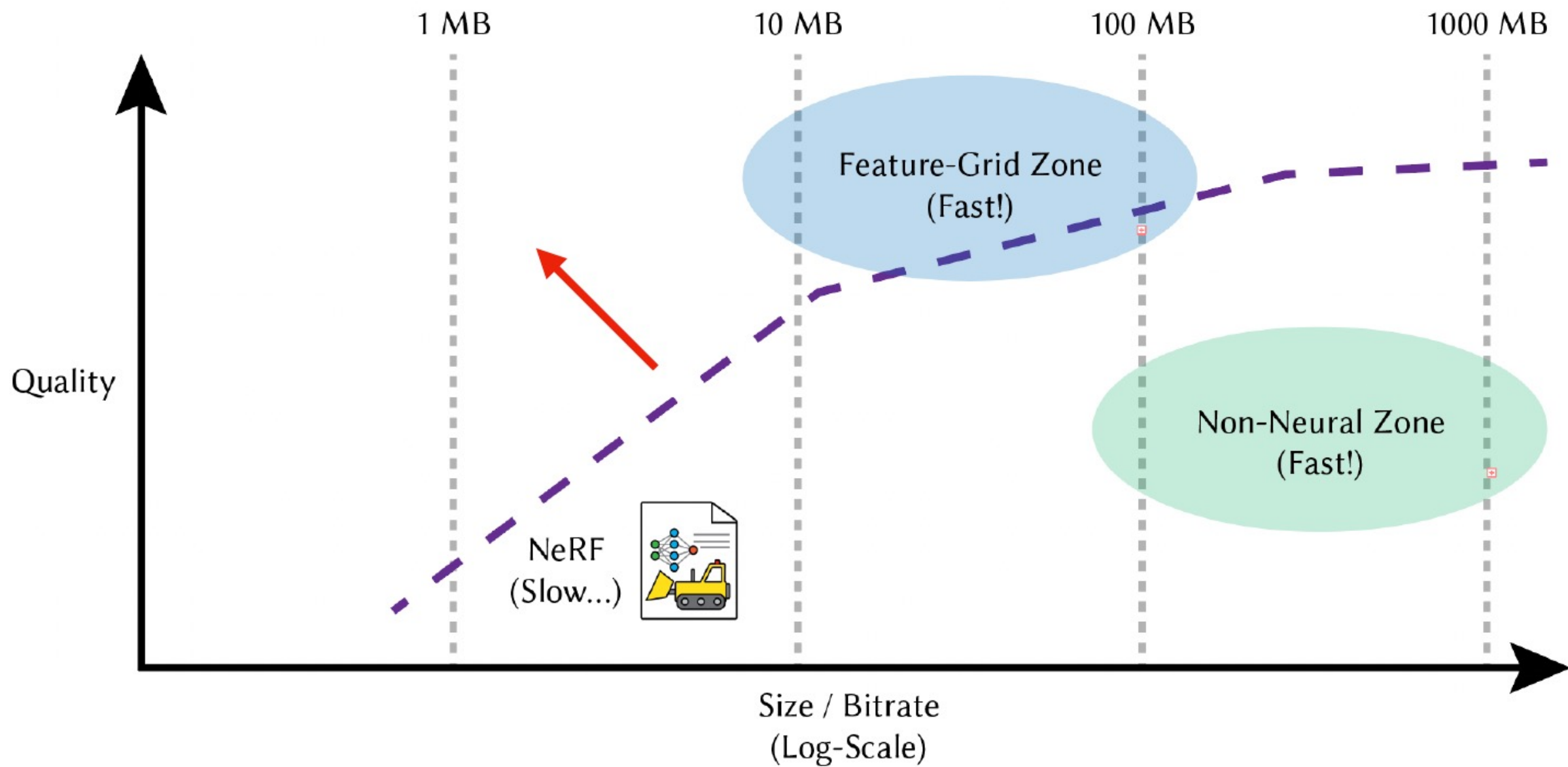
# Hash Grids



Codebook

[x y z]

Hash

[Instant-NGP (Muller et al.)]

Pros:

- Densely supported
- Disaggregate resolution from memory cost
- No complex data structures
- Performant memory access if codebook is small enough

Cons:

- Multiresolution and large codebooks needed for collision resolution
- Features not spatially local

# Key limitations of the original NeRF

- Very slow in training and inference
- Requires Ground-Truth poses
- Do not generalize to new scenes

**BARF 🤮: Bundle-Adjusting Neural Radiance Fields**

Chen-Hsuan Lin😷   Wei-Chiu Ma🤢   Antonio Torralba🤢   Simon Lucey😷🤔

😷Carnegie Mellon University   🤢Massachusetts Institute of Technology   🤔The University of Adelaide

IEEE International Conference on Computer Vision (**ICCV**), 2021
**oral presentation**

# 🌰 NeRF in a nutshell
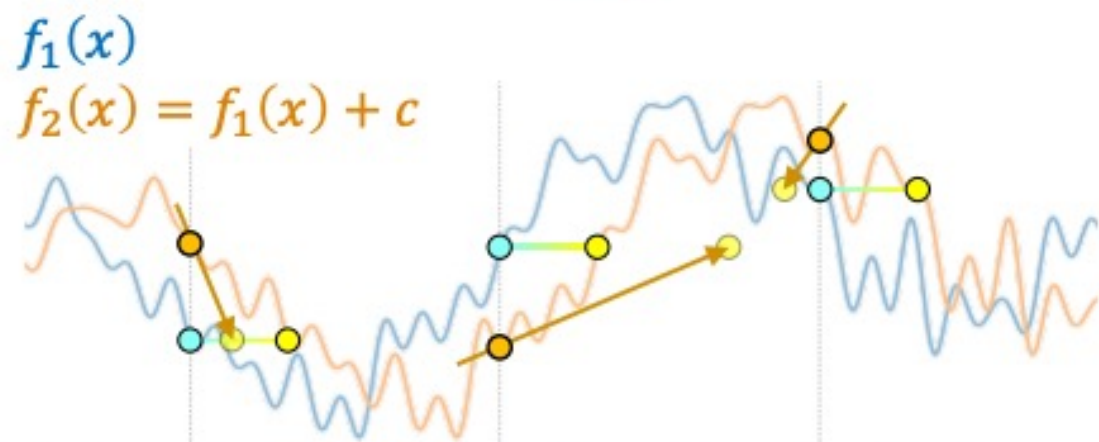
$$\min_{\mathbf{p}_1,\dots,\mathbf{p}_M,\boldsymbol{\Theta}} \sum_{i=1}^{M} \sum_{\mathbf{u}} \left\| \hat{\mathcal{I}}(\mathbf{u};\mathbf{p}_i,\boldsymbol{\Theta}) - \mathcal{I}_i(\mathbf{u}) \right\|_2^2$$

frames $M$

RGB (rendered)  camera pose  network params.  RGB

We want to optimize the poses as well!



image plane

$\mathbf{u}$

camera center

input 3D coordinates $\mathbf{x}$ $\gamma(\mathbf{x})$

neural network (MLP)

$\mathbf{c}$ color (RGB)

$\boldsymbol{\sigma}$ volume density

**positional encoding**

✅ encourages representation with high frequency

❌ detrimental to gradient-based **registration**!!

$f_1(x)$

$f_2(x) = f_1(x) + c$

❌ gets stuck in suboptimal solutions

✅ smooth signals ➔ coherent updates

V 2020

**SOLUTION** 🤩 :
make it **coarse-to-fine**!

Resolve large pose misalignment & coarse scene representation

Gradually activate higher-frequency components in positional encoding

Refine granular pose misalignment & high-fidelity scene representation

# Key limitations of the original NeRF

- Very slow in training and inference
- Requires Ground-Truth poses
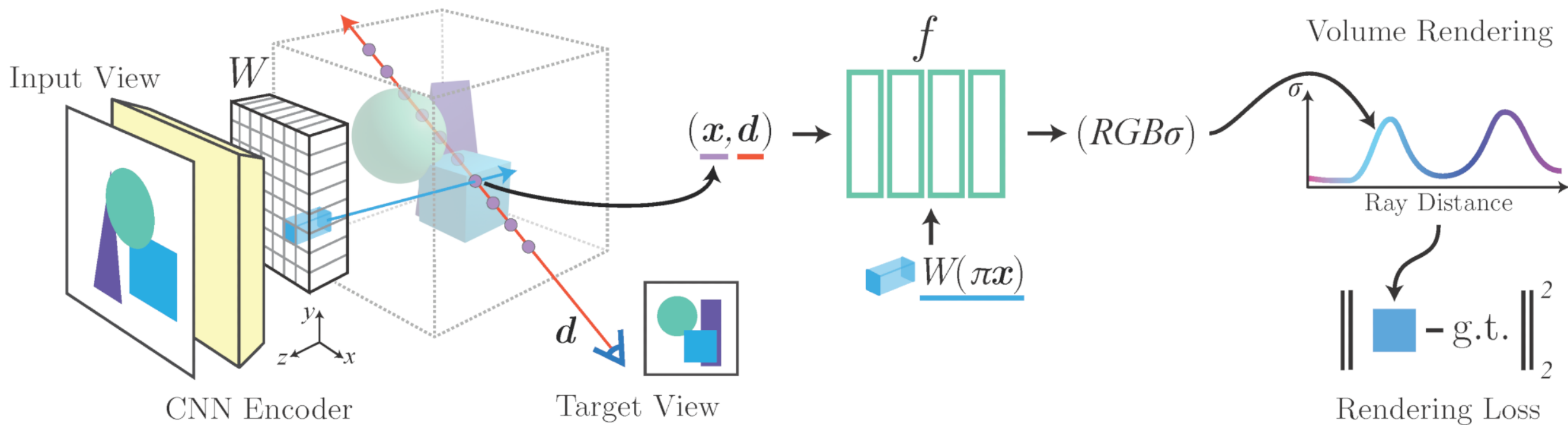- Do not generalize to new scenes

# pixelNeRF

## Neural Radiance Fields from One or Few Images

CVPR 2021

Alex Yu    Vickie Ye    Matthew Tancik    Angjoo Kanazawa
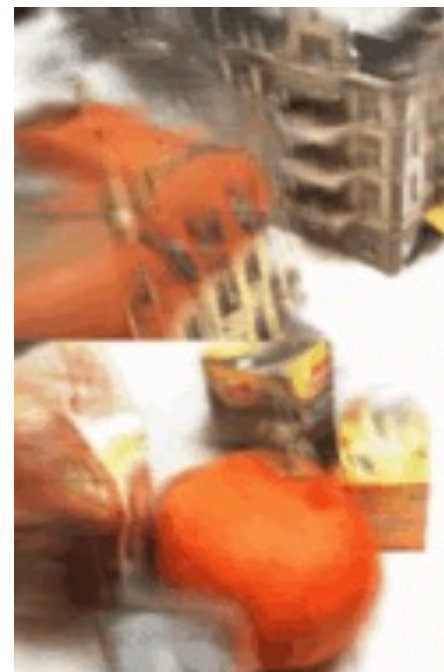
UC Berkeley

Input Images       PixelNeRF       NeRF

# Slide Credits

- "Introduction to Computer Vision", Noah Snavely, Cornell Tech, Spring 2022

- "Understanding and Extending Neural Radiance Field", Jon Barron MIT & Tu Munich Lecture.

- "Neural Fields in Computer Vision", CVRP 2022 Tutorial.

- Shubham Tulsiani, "Learning for 3D Vision", Spring 2022, CMU

- Leo Guibas, JJ Park, "Neural Models for 3D geometry", Spring 2022, Stanford.