

Lecture 24: Neural Radiance Fields (NeRFs)

COMP 590/776: Computer Vision

Instructor: Roni Sengupta

Final-exam focus

Schedule

Date	Topic	Details	Special Dates
Intro & Review			
Tues Aug 22	Intro. to Computer Vision & Ethical Concerns	Lecture Slide	
Thrs Aug 24	Maths Review (Linear Algebra, Probability, Calculus)	Lecture Slide [pdf]	HW1: Maths review (assigned)
Colors & Imaging			
Tues Aug 29	No Class		
Thrs Aug 31	Color & Color Spaces	Lecture Slide [pdf]	
Tues Sept 5	No Class (Well-being day)		HW1: Maths review (due)
Thrs Sept 7	In-Camera Imaging Pipeline	Lecture Slide [pdf]	
Image Processing			
Tues Sept 12	Filtering - Convolution, Gradients, & Edges	Lecture Slide [pdf]	
Thrs Sept 14	Frequency domain - Fourier Analysis	Lecture Slide [pdf]	
Features			
Tues Sept 19	Feature Detection (Corner & Blob)	Lecture Slide [pdf]	
Thrs Sept 21	Feature Descriptor & Matching (SIFT)	Lecture Slide [pdf]	

- Convolution Operator
- Different Filters
- Aliasing
- Image Derivatives
- Fourier Transform
- Canny Edge Detection
- Harris Corner Detector and it's properties

Final-exam focus

		2D Transformation	
Tues Sept 26	2D Transformations & Fitting	Lecture Slide [pdf]	
Thrs Sept 28	RANSAC + Image Blending	Lecture Slide [pdf]	
		Learning & Perception	
Tues Oct 3	Recognition	Lecture Slide [pdf]	
Thrs Oct 5	Detection	Lecture Slide [pdf]	
Tues Oct 10	Segmentation & Matting	Lecture Slide [pdf]	
Thrs Oct 12	No Class (University Day)		HW3: Panorama (due Friday Oct 13) HW4: Deep Learning (assigned)
Tues Oct 17	Generative Models in Computer Vision	Lecture Slide [pdf]	Online Lecture
Thrs Oct 19	No Class (Fall Break)		

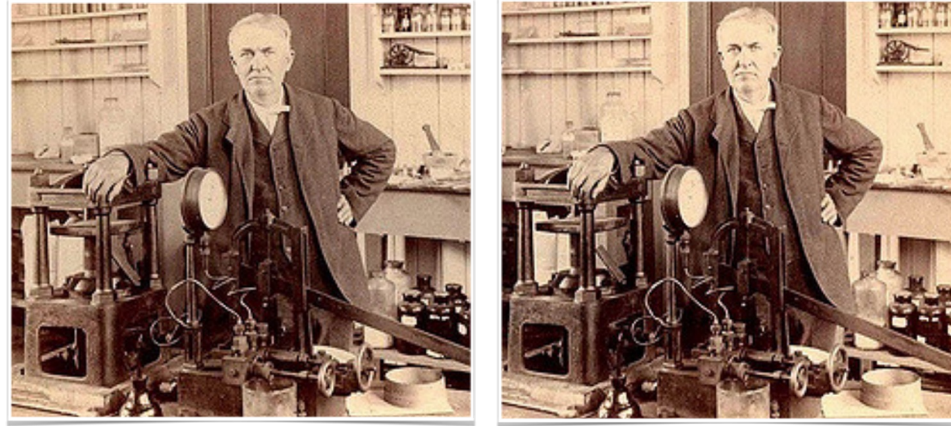
- 2D Transformation: Affine & Homography
- Least Square fitting
- Estimating Homography
- Basic Structures of a CNN
- How to train a CNN for classification
- Different kinds of object detectors (high-level)

Final-exam focus

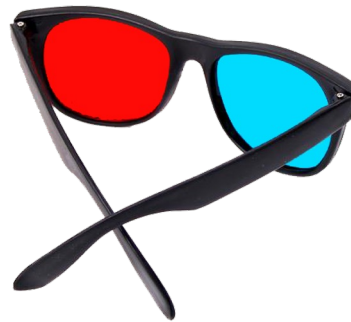
Tues Oct 24	Camera Models + Calibration - 1		
Thrs Oct 26	Camera Models + Calibration - 2		
Tues Oct 31	Two-view Geometry-1		
Thrs Nov 2	Two-view Geometry-2		
Tues Nov 7	Stereo		
Thrs Nov 9	Multi-view Stereo		
Tues Nov 14	Structure from Motion		
Thrs Nov 16	Light & Photometric Stereo		
Tues Nov 21	Deep Learning for MVS, SfM, PS		
Thrs Nov 23	No Class (Thanksgiving)		
Tues Nov 28	NeRFs		HW6: 3D Vision 2 (due)
Thrs Nov 30	Mid-term Review		
Tues Dec 5	Final Exam (in-class and/or take home)		Syllabus: Whole course

- World->camera->image coord.
- Perspective distortion
- Vanishing points, lines and planes
- Essential & Fundamental Matrix and its properties
- Normalized 8 point algorithm
- Depth-disparity relation (stereo)
- Image rectification (stereo)
- SfM: ambiguities and minimal view-points for solving.
- High-level understanding of different 3D reconstruction techniques and their relative advantage, disadvantage, and used cases.

Stereo Photography



Viewing Devices



Left



Right





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

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

NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis

ECCV 2020



Ben Mildenhall*



UC Berkeley



Pratul Srinivasan*



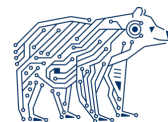
UC Berkeley



Matt Tancik*



UC Berkeley



Jon Barron



Google Research



Ravi Ramamoorthi



UC San Diego



Ren Ng



UC Berkeley



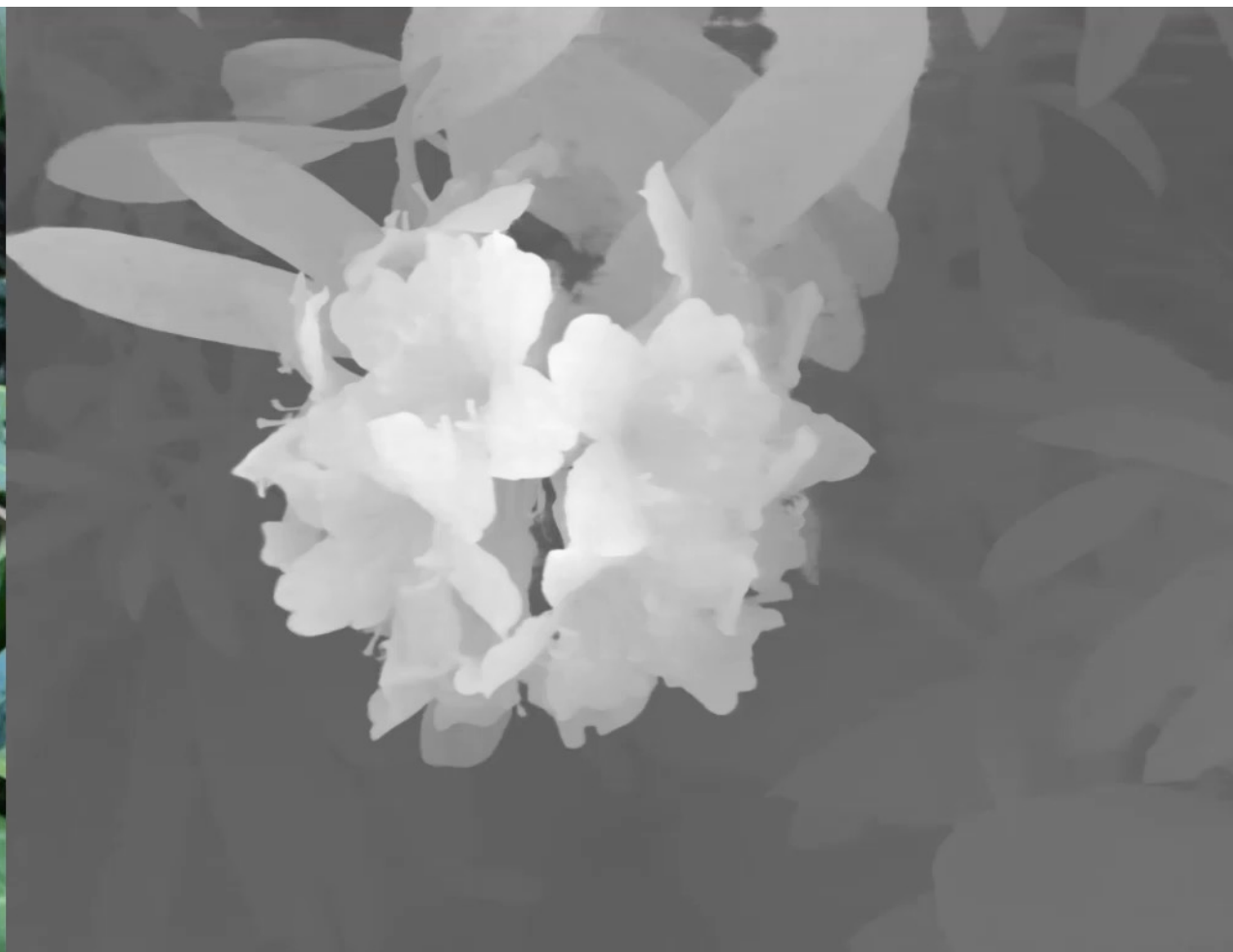


Given a set of sparse views of an object with known camera poses

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.



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.

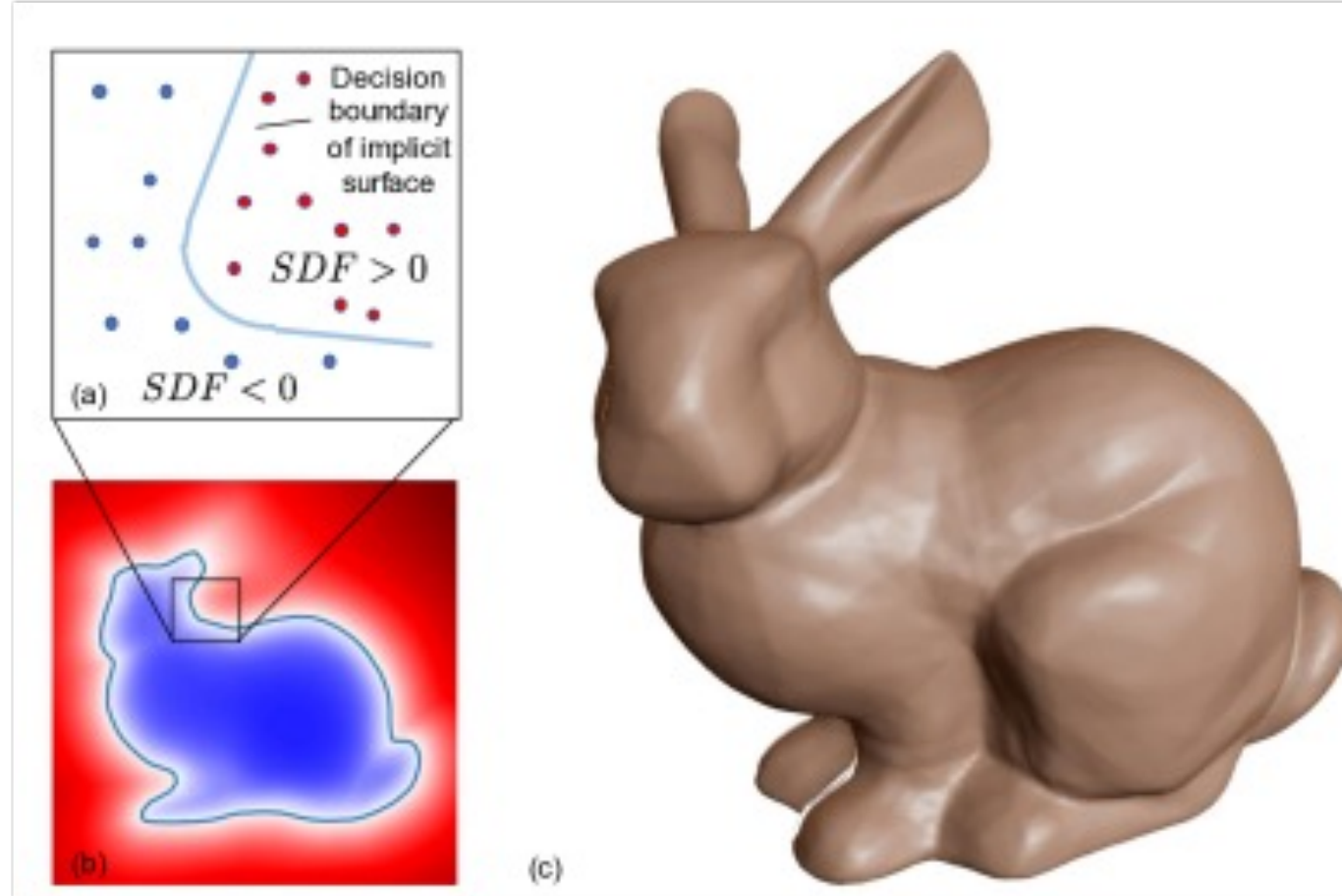


Neural 3D Video Synthesis
from Multi-view Video,
Li et al., CVPR 2022

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

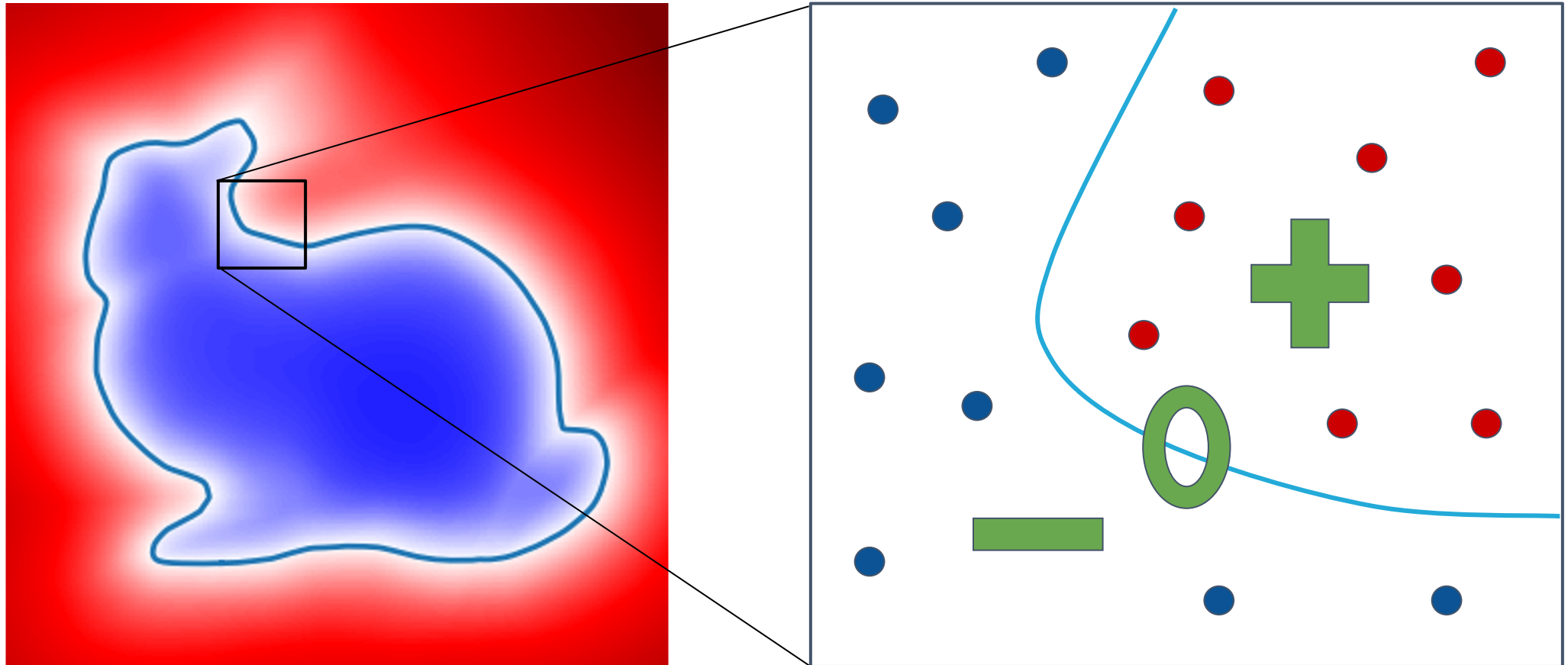
$SDF(X) = 0$, when X is on the surface.
 $SDF(X) > 0$, when X is outside the surface
 $SDF(X) < 0$, when X is inside the surface

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

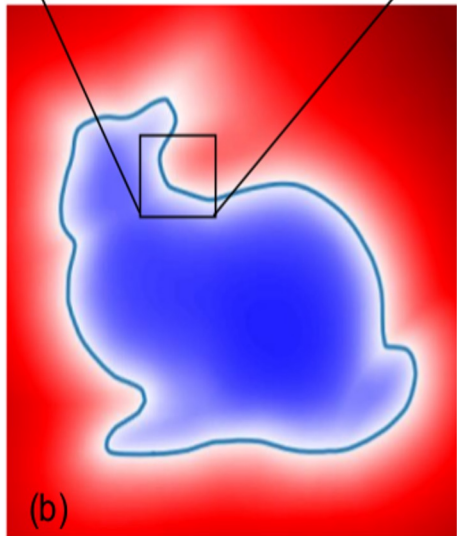
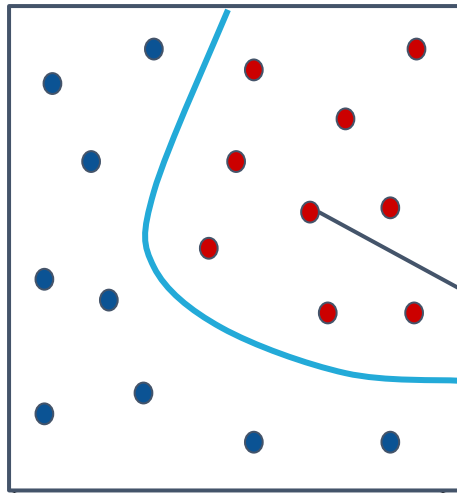


Deep SDF: Use a neural network (co-ordinate based MLP) to represent the SDF function.

Signed Distance Function



Regression of Continuous SDF

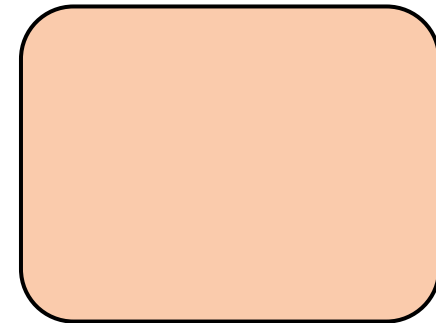


(b)

(x, y, z)



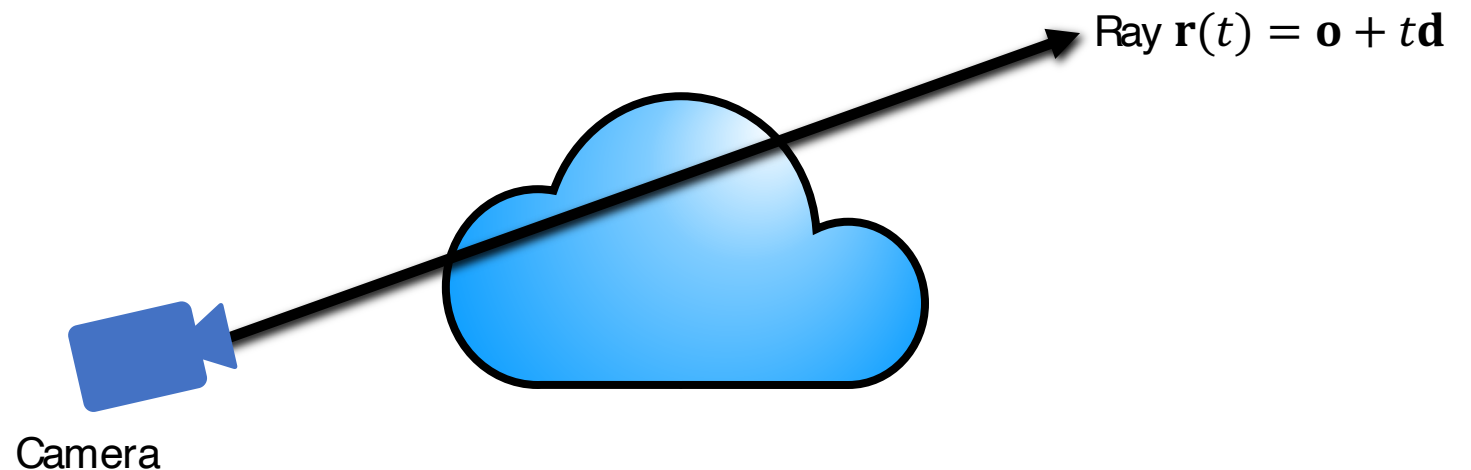
NN



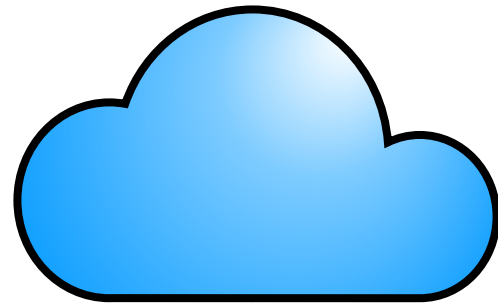
SDF

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?

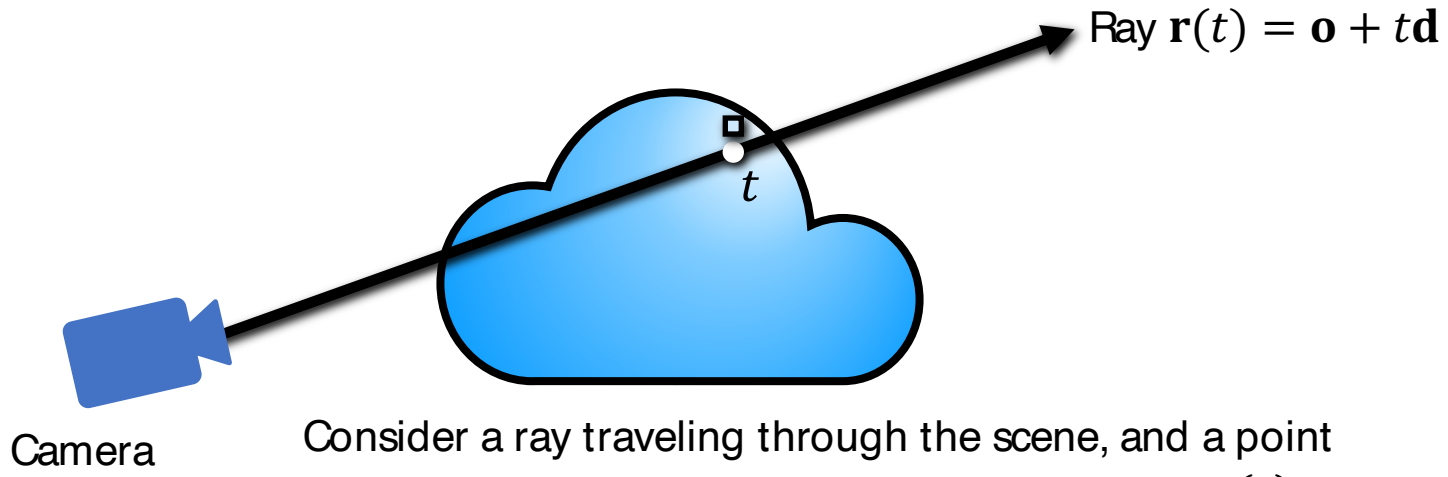


Volumetric formulation for NeRF



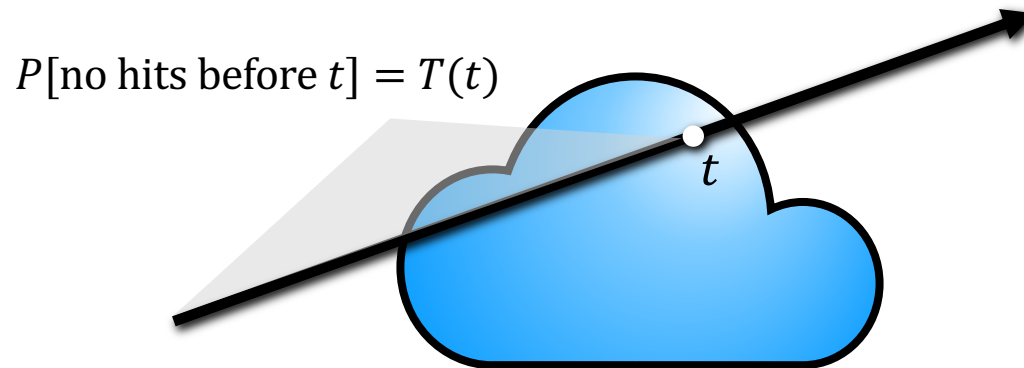
Scene is a cloud of colored fog

Volumetric formulation for NeRF



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) $\hat{L}(t)$

Volumetric formulation for NeRF



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

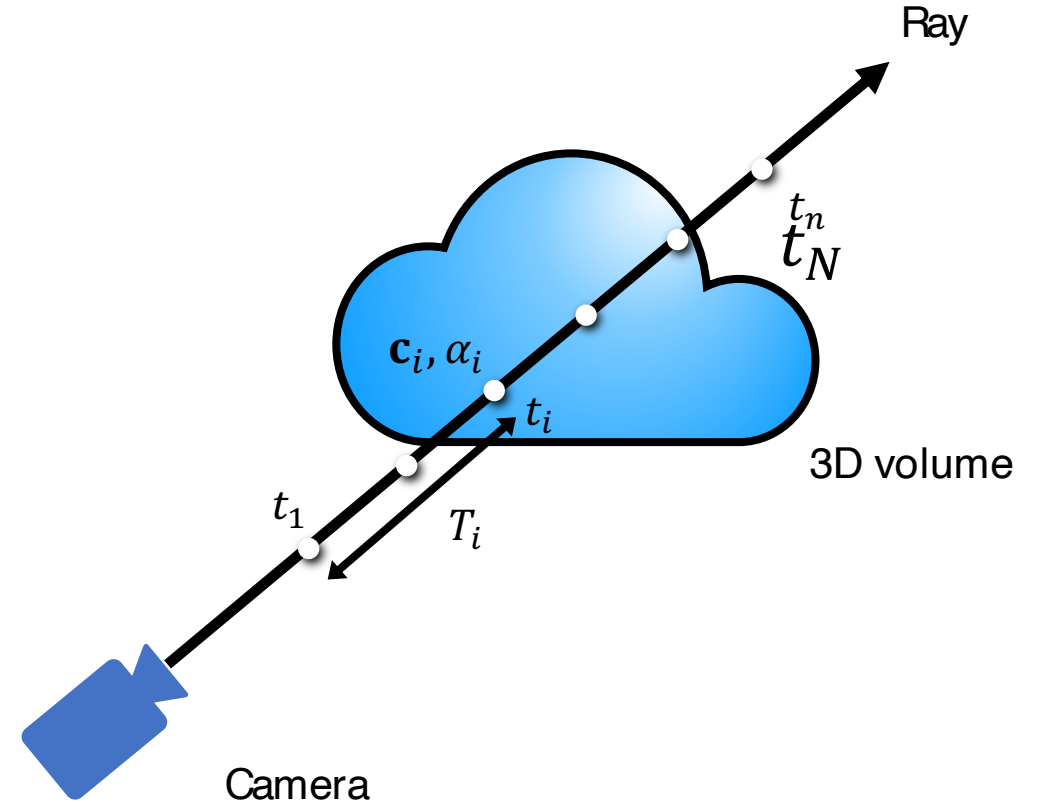
$$\mathbf{c} \approx \sum_{i=1}^n T_i \alpha_i \mathbf{c}_i$$

final rendered color along ray \mathbf{c}
 weights T_i
 colors \mathbf{c}_i

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



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

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

Volume rendering estimation: integrating color along a ray

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

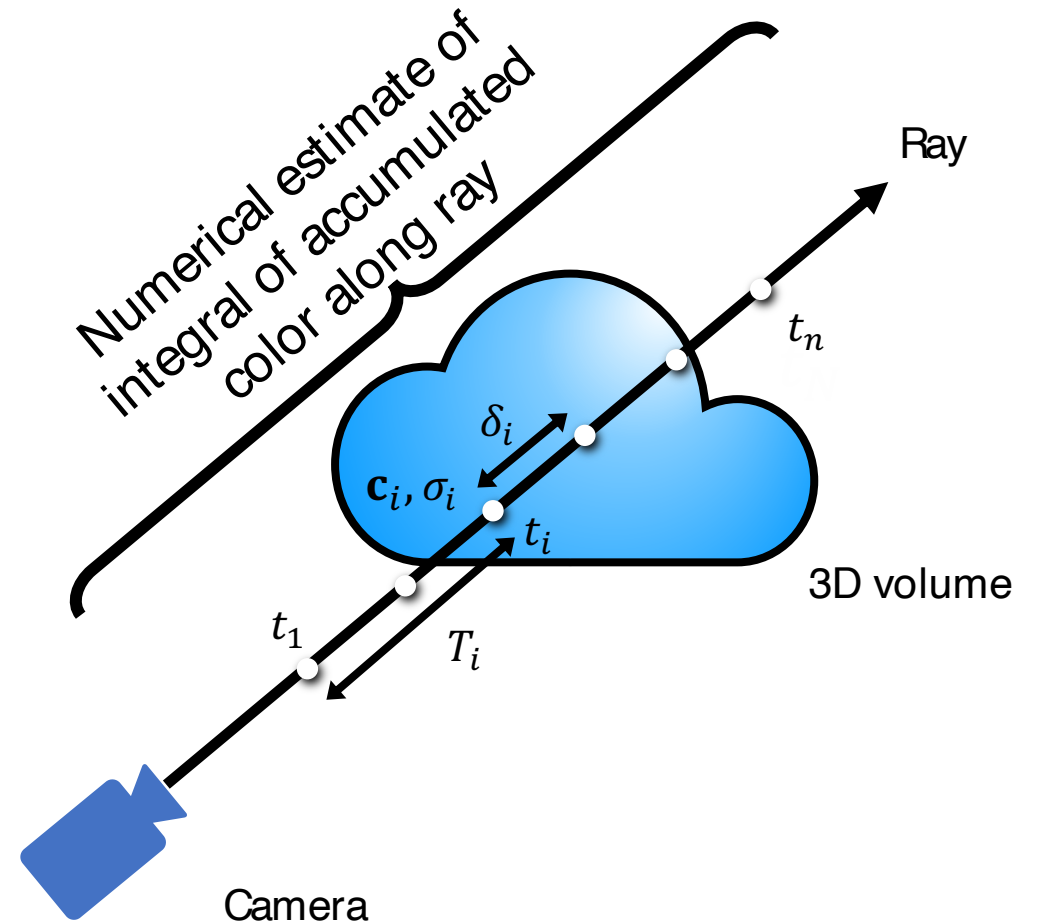
$$\mathbf{c} \approx \sum_{i=1}^n T_i \alpha_i \mathbf{c}_i$$

final rendered color along ray weights colors

How much light is blocked earlier along ray:

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

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



Volume rendering estimation: integrating color along a ray

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

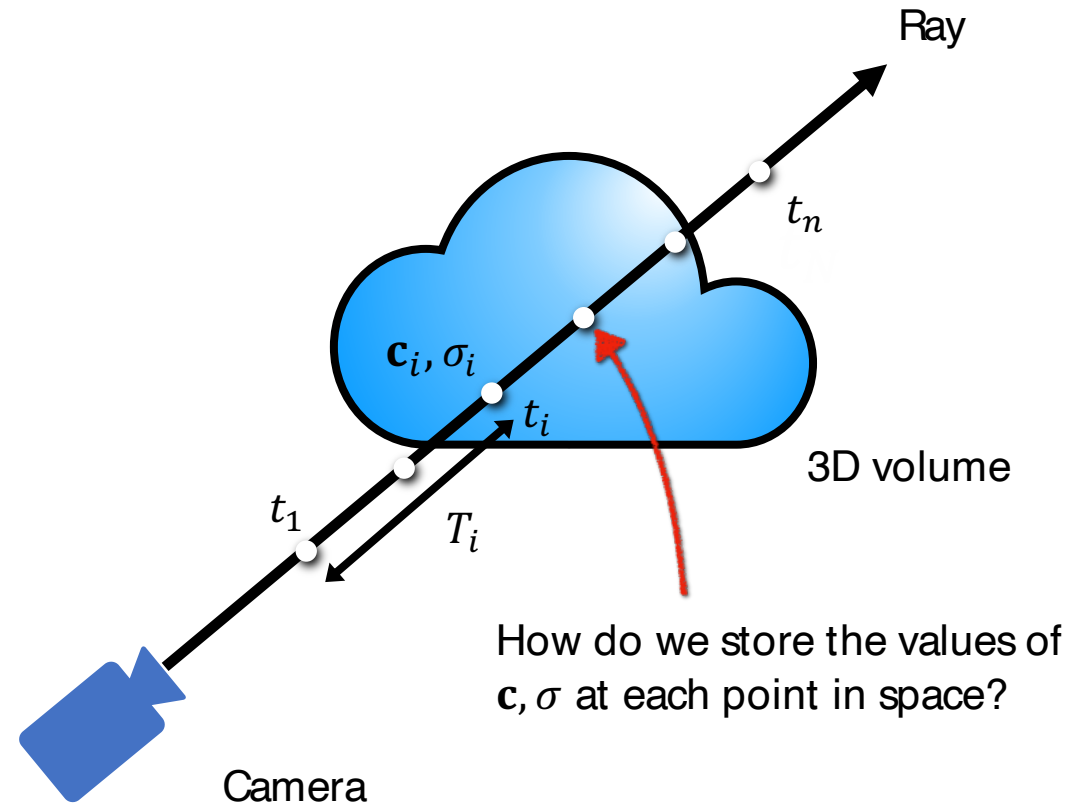
$$\mathbf{c} \approx \sum_{i=1}^n T_i \alpha_i \mathbf{c}_i$$

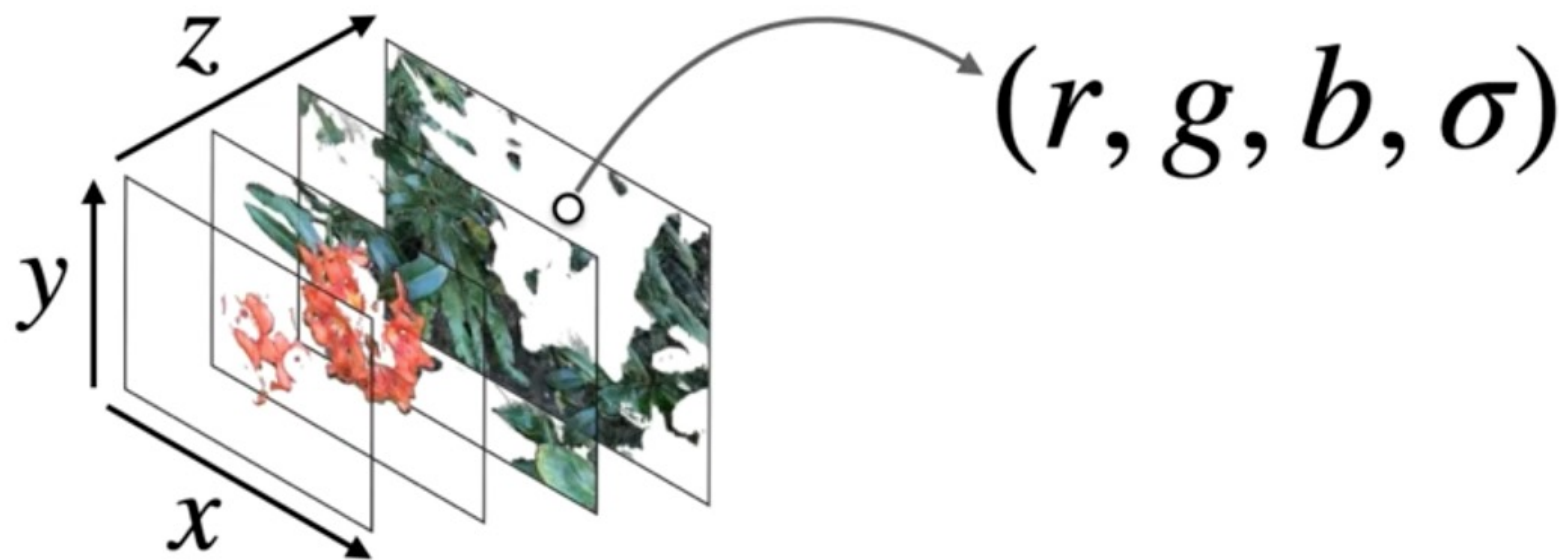
final rendered color along ray weights colors

How much light is blocked earlier along ray:

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

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



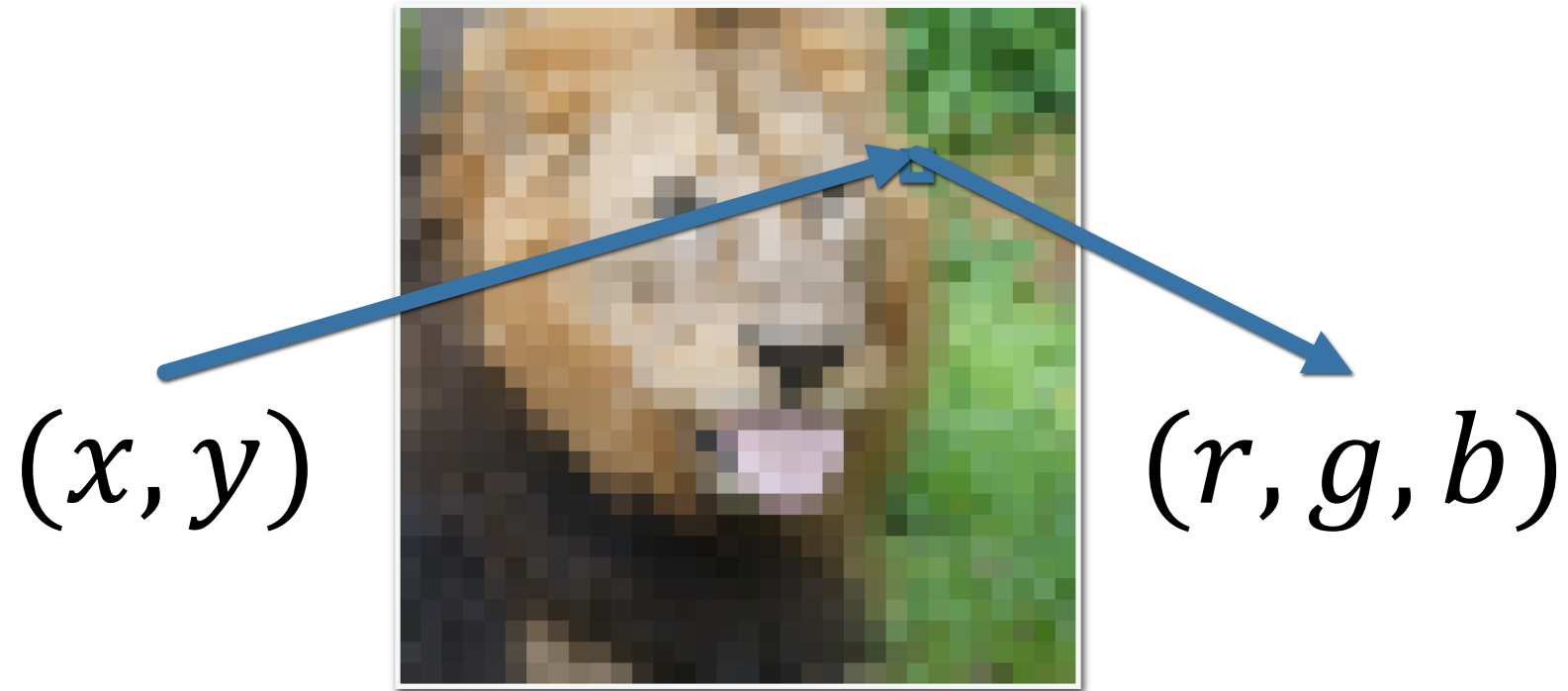


versus

$$(x, y, z, \theta, \phi) \rightarrow \text{[blue bars]} \rightarrow (r, g, b, \sigma)$$

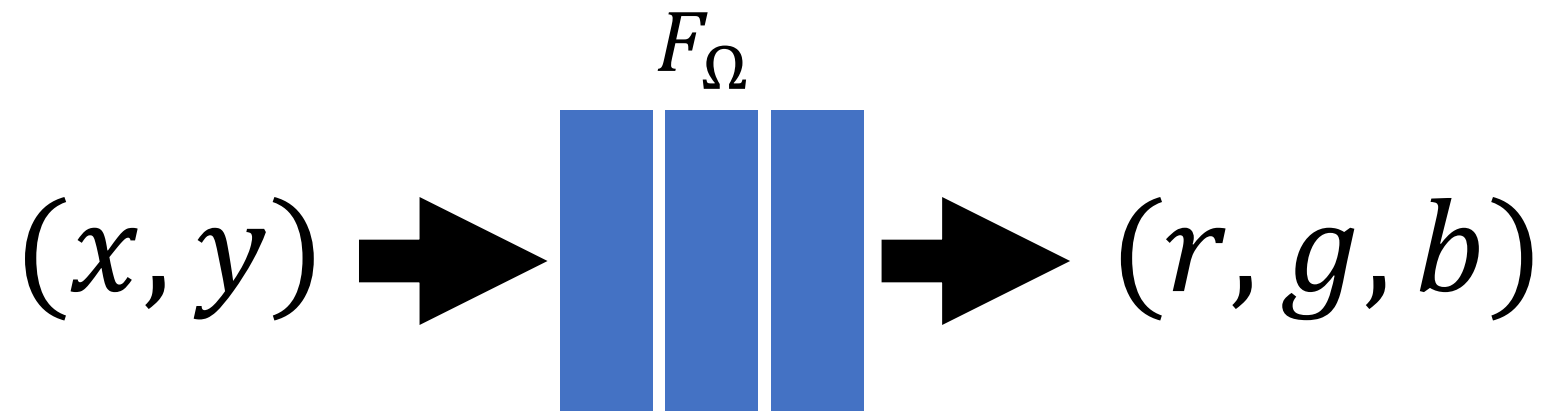
F_{Ω}

Toy problem: storing 2D image data



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

Toy problem: storing 2D image data



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

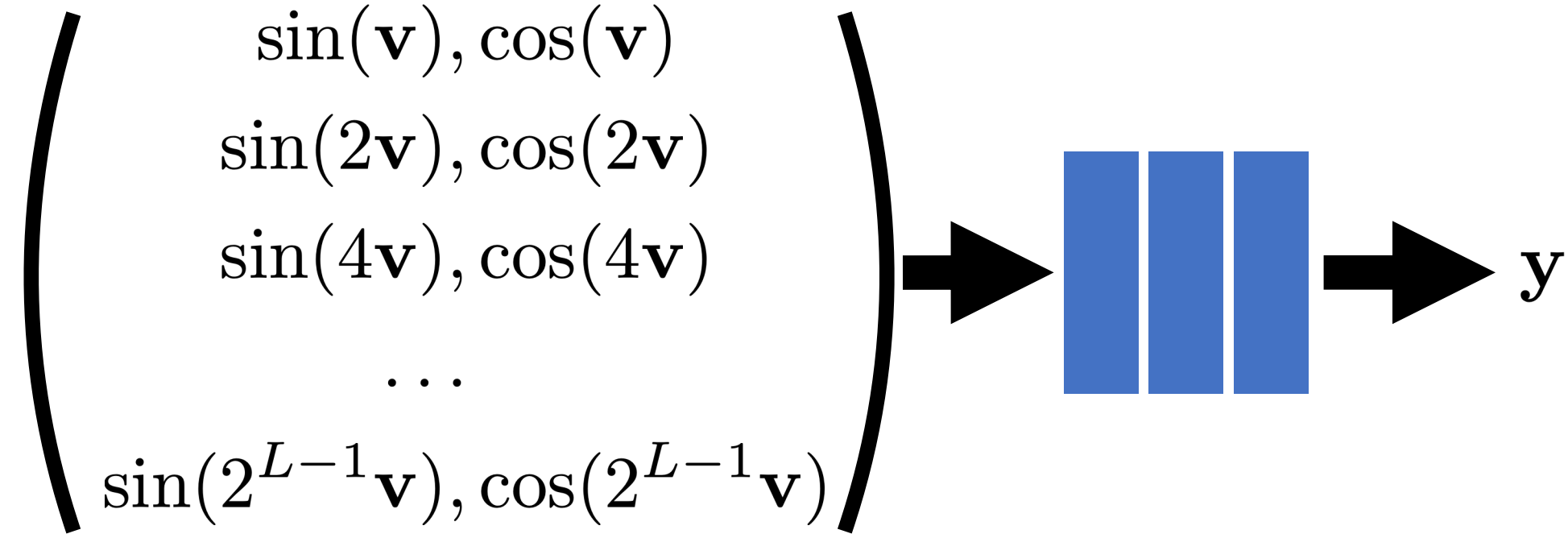
Problem:

- “Standard” coordinate-based MLPs cannot represent high frequency functions.

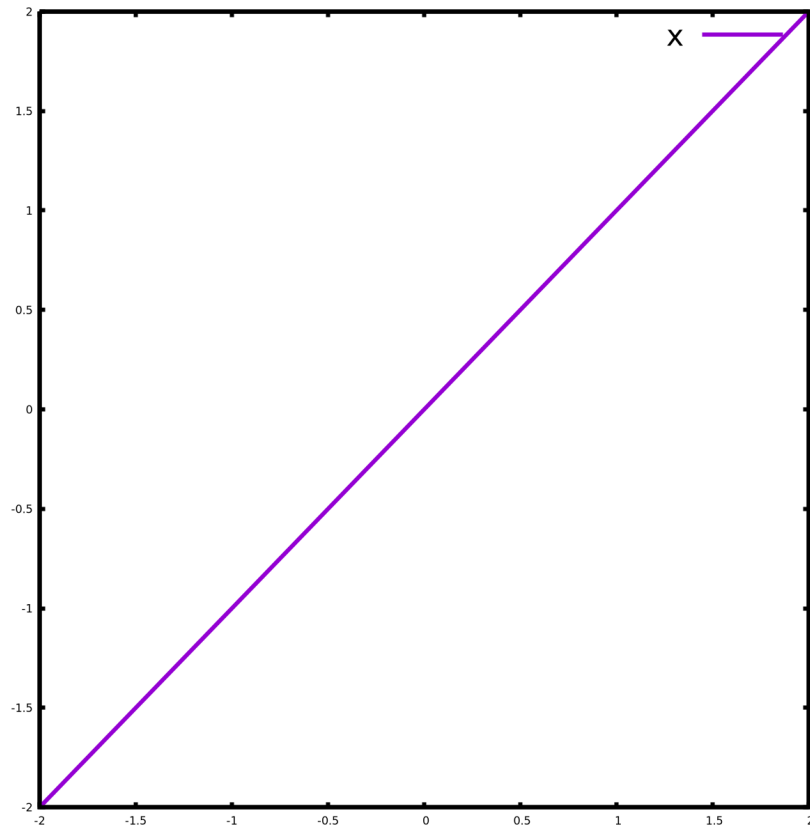
Solution:

- Pass input coordinates through a high frequency mapping first.

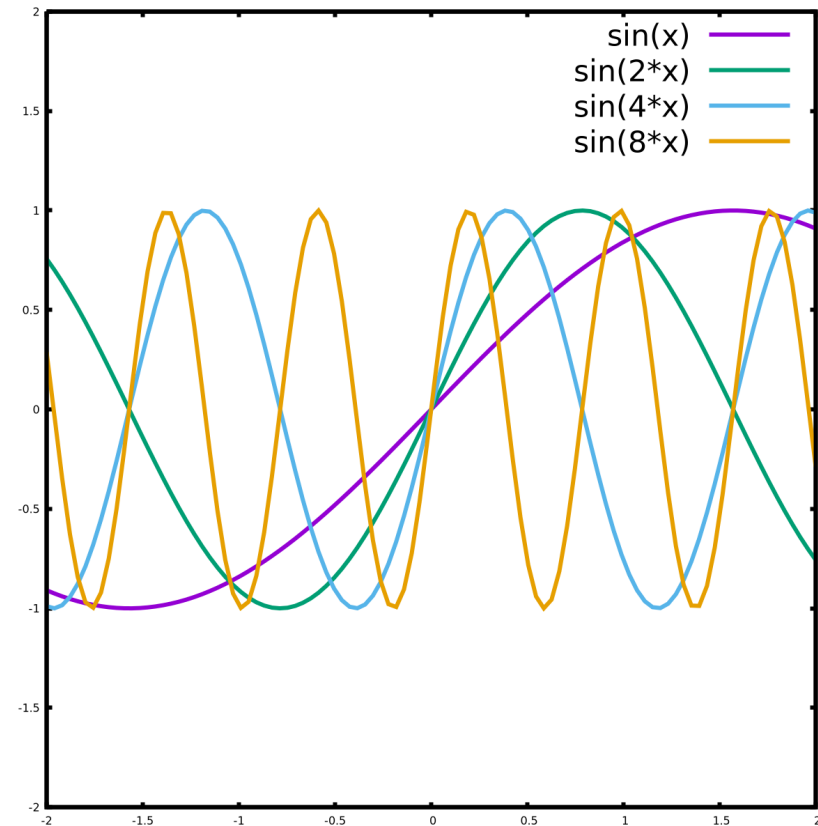
Example mapping: "positional encoding"



Positional encoding



Raw encoding of a number x



“Positional encoding” of a number x

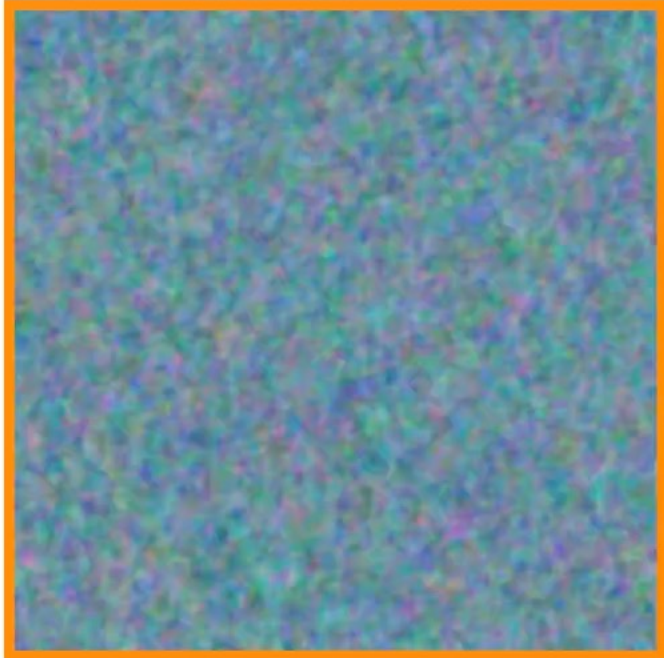
Problem solved!



Ground truth image

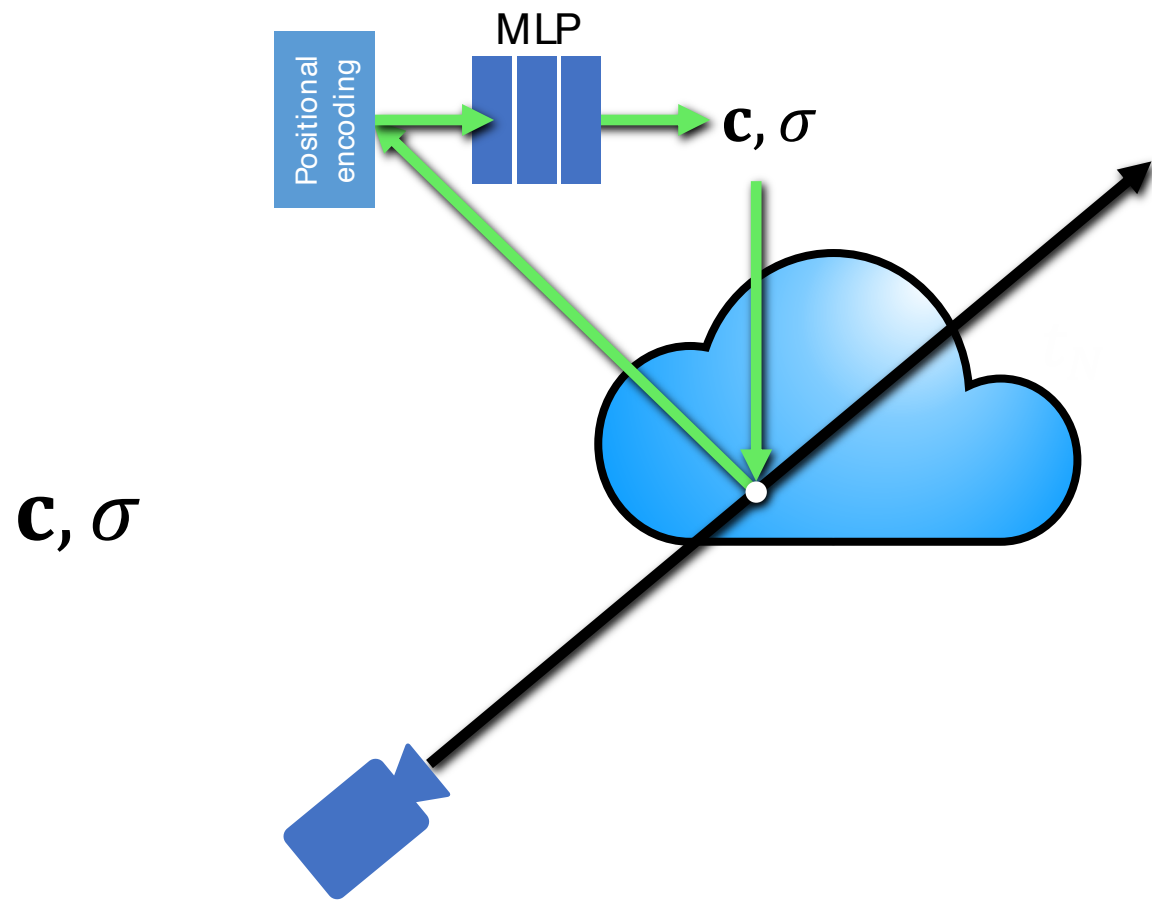


Neural network output without high frequency mapping

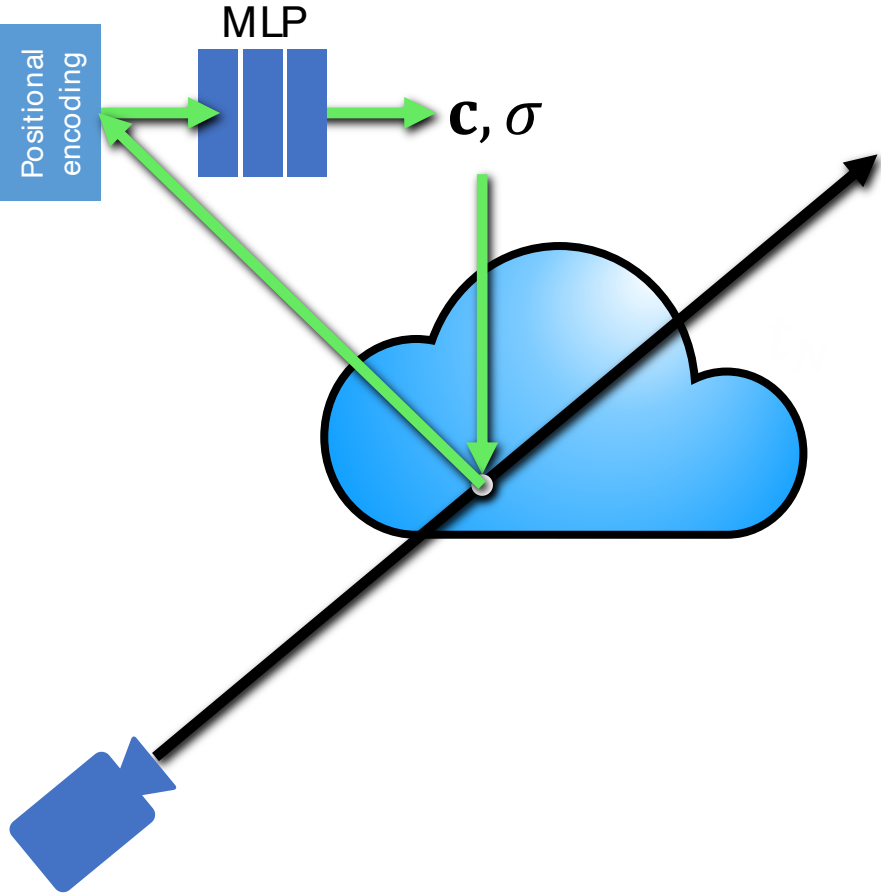


Neural network output with high frequency mapping

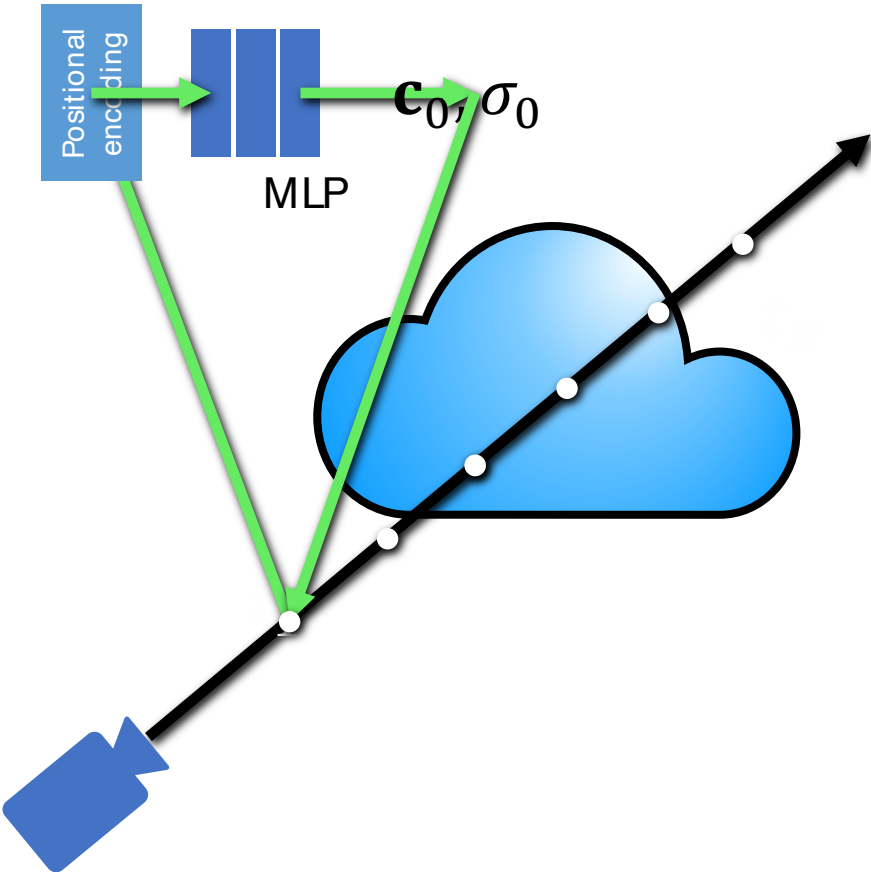
NeRF = volume rendering +
coordinate-based network



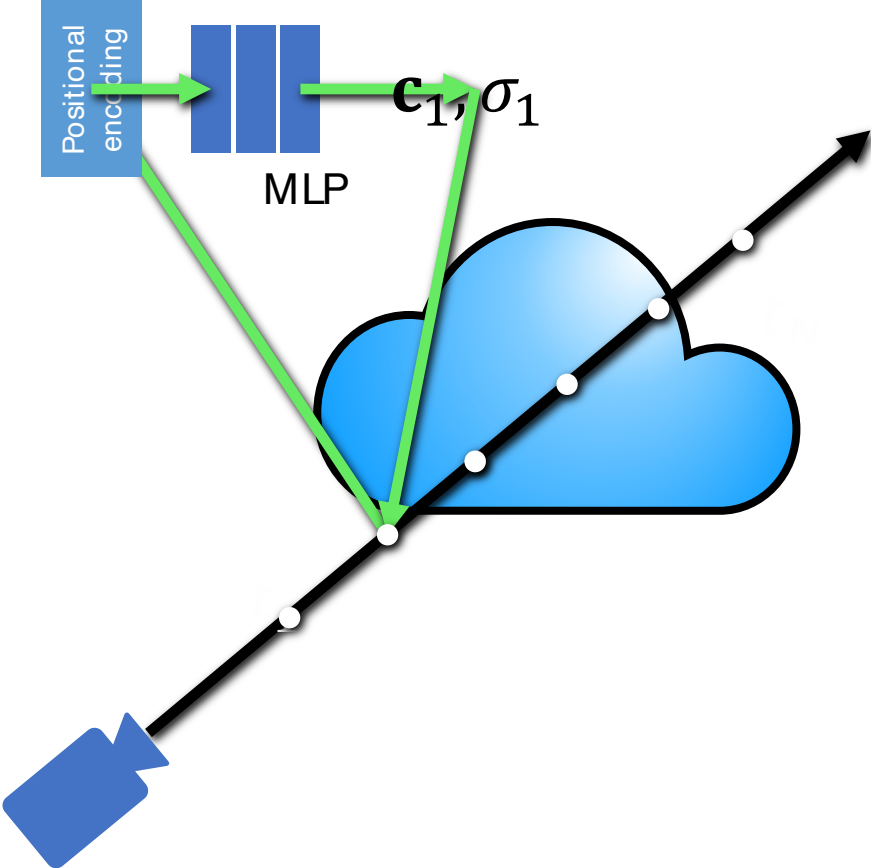
\mathbf{c}, σ



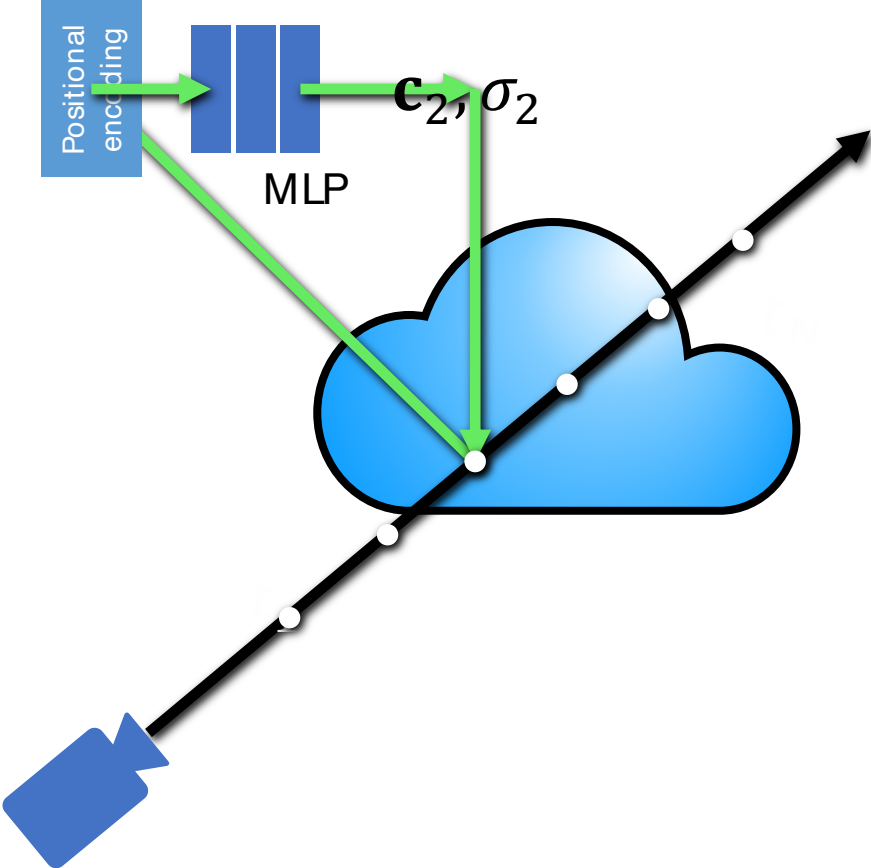
\mathbf{c}, σ



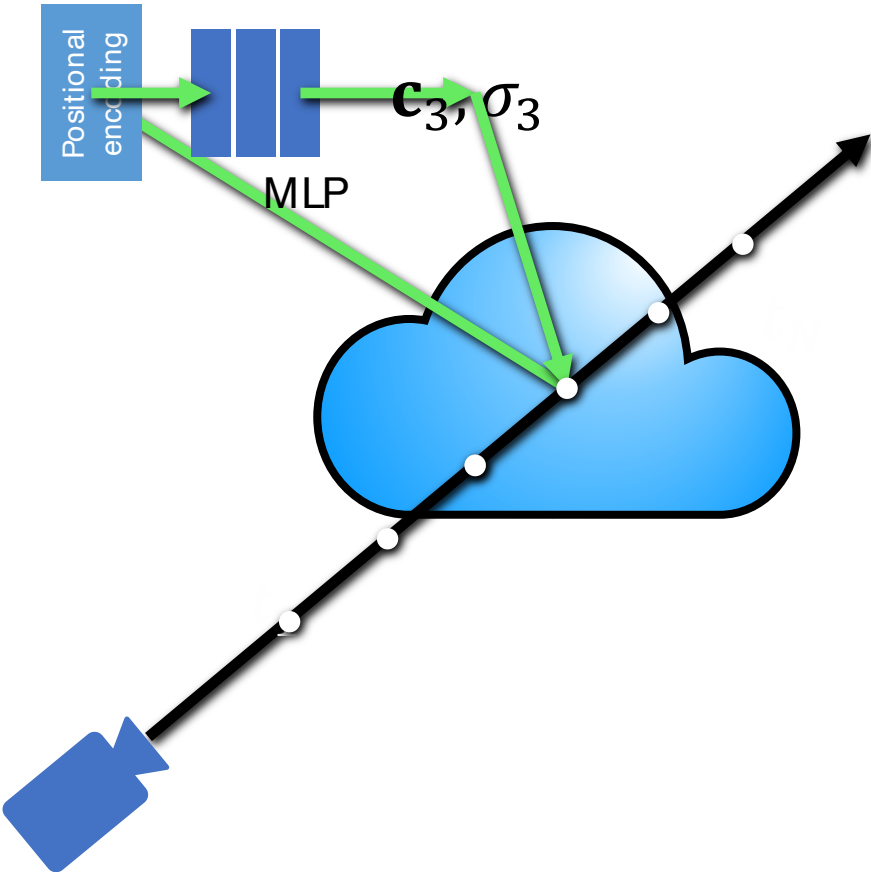
\mathbf{c}, σ



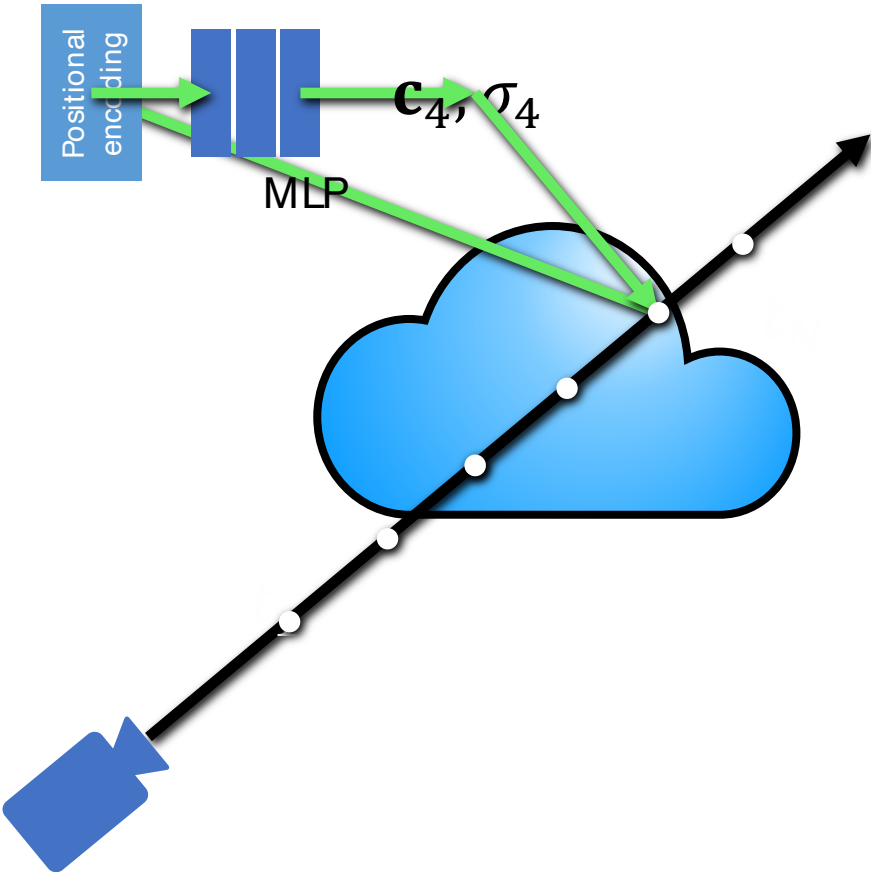
\mathbf{c}, σ



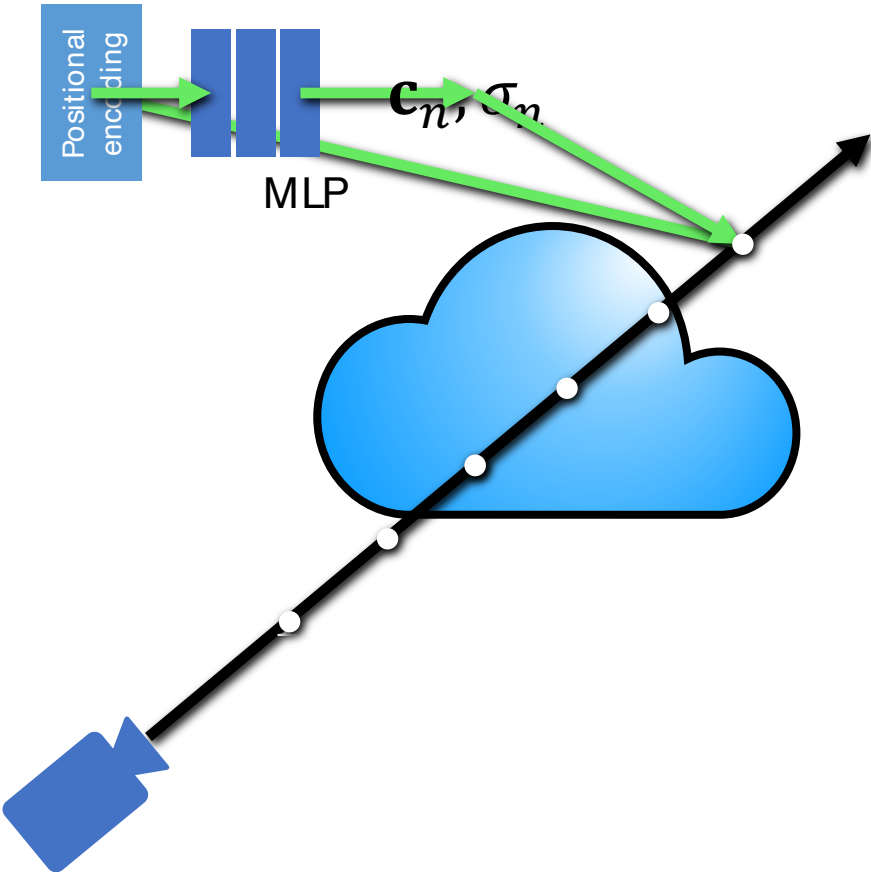
\mathbf{c}, σ

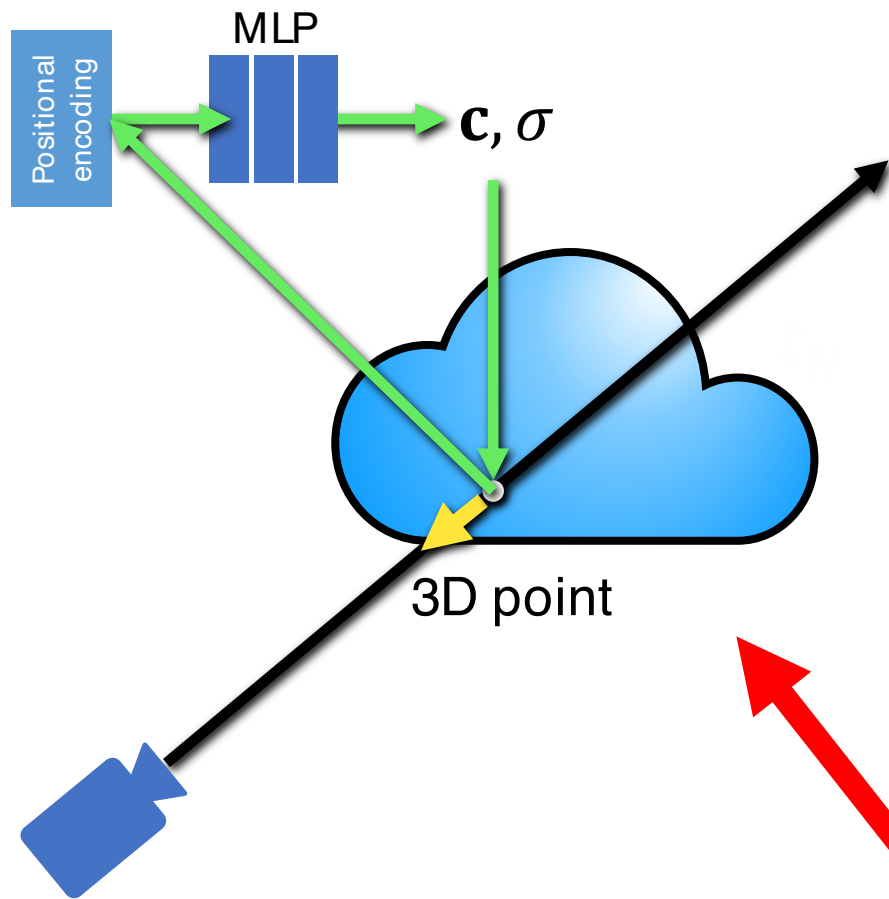


\mathbf{c}, σ



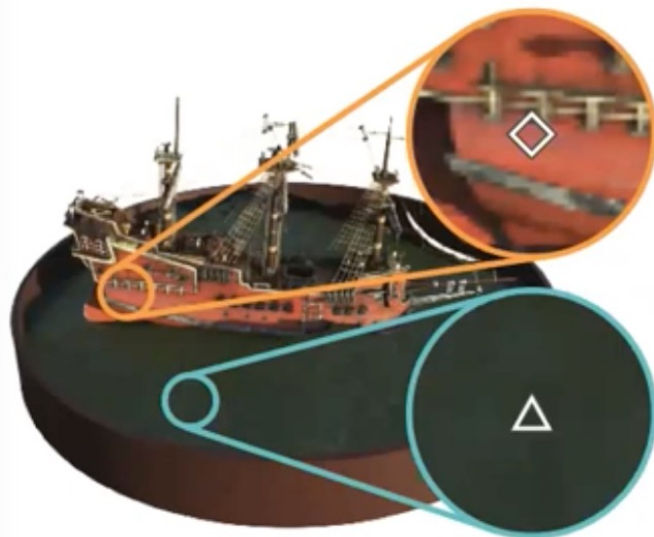
\mathbf{c}, σ



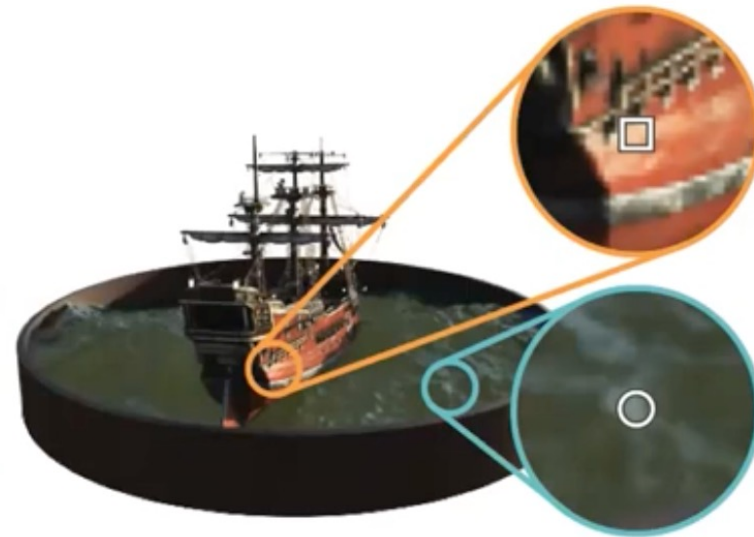


Include the ray direction in the input to the MLP → allows for capturing and rendering view-dependent effects (e.g., shiny surfaces)

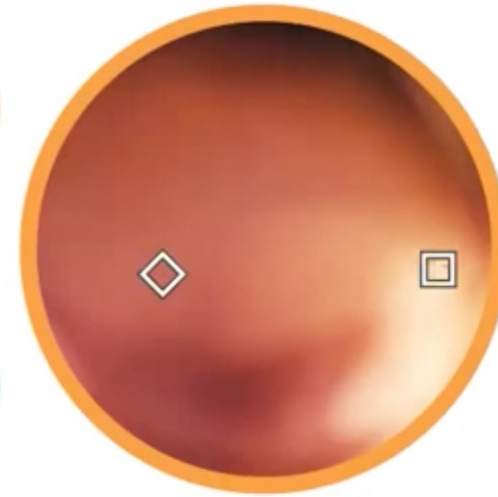
Modeling view dependent effects



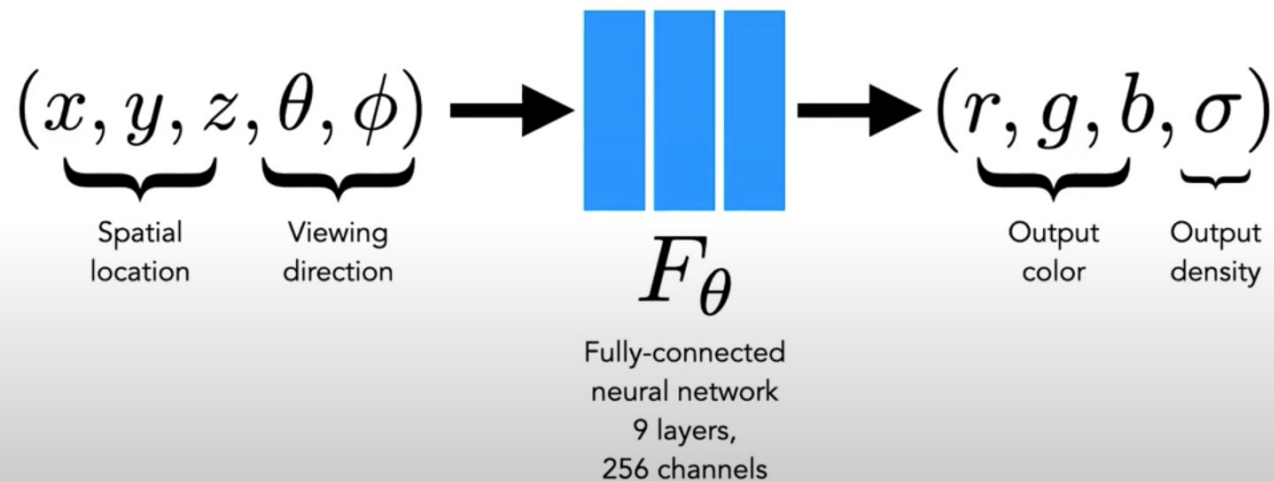
(a) View 1



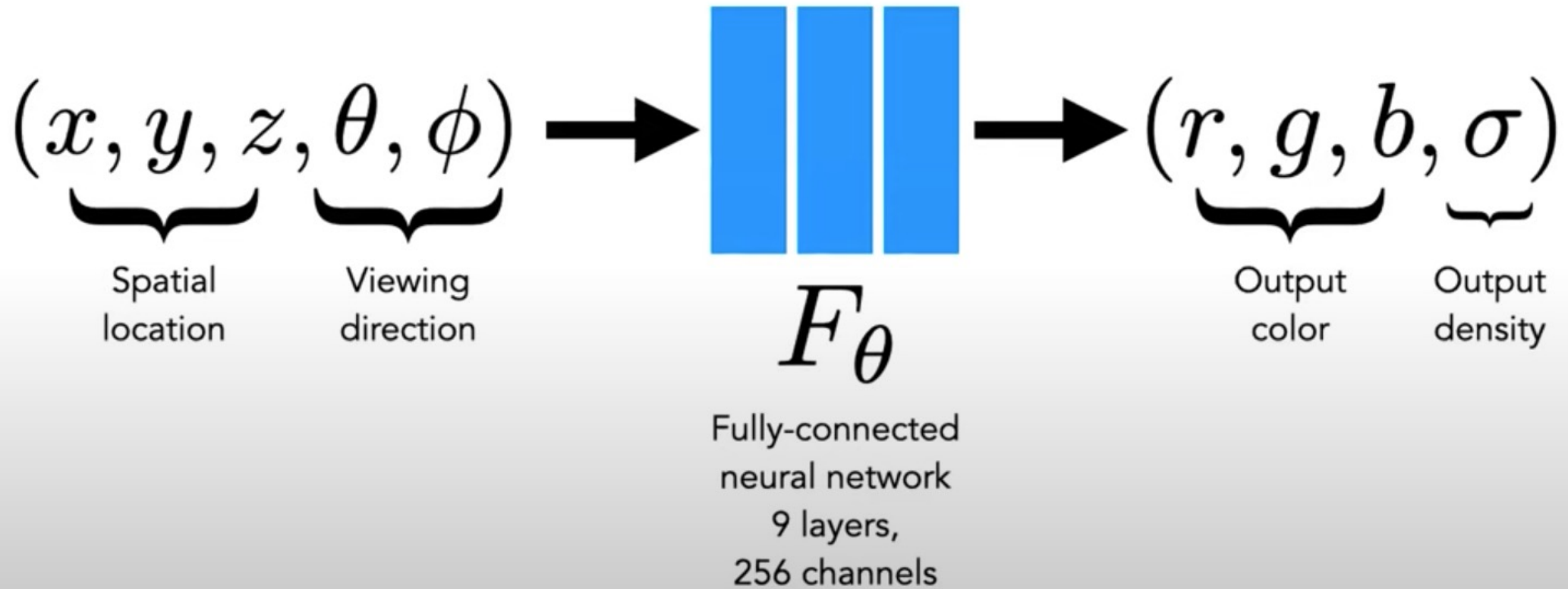
(b) View 2



(c) Radiance Distributions

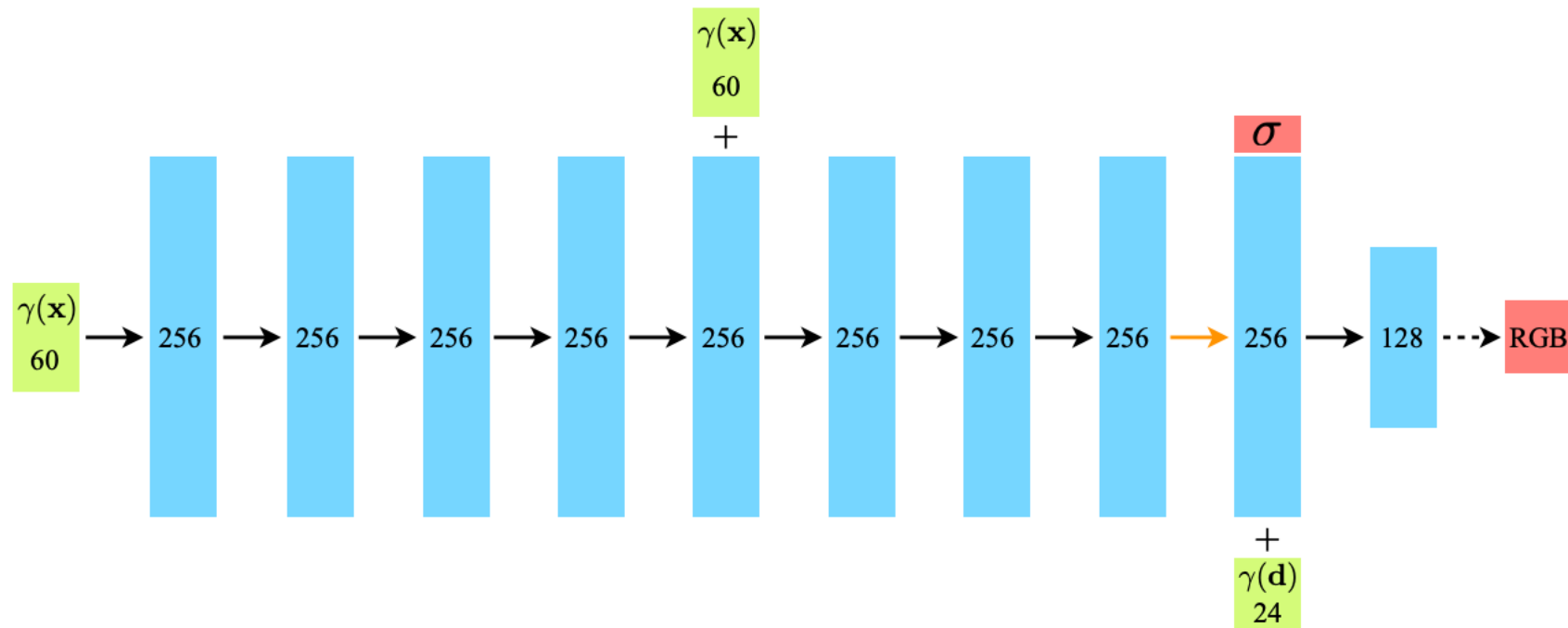


What do we learn in NeRF?

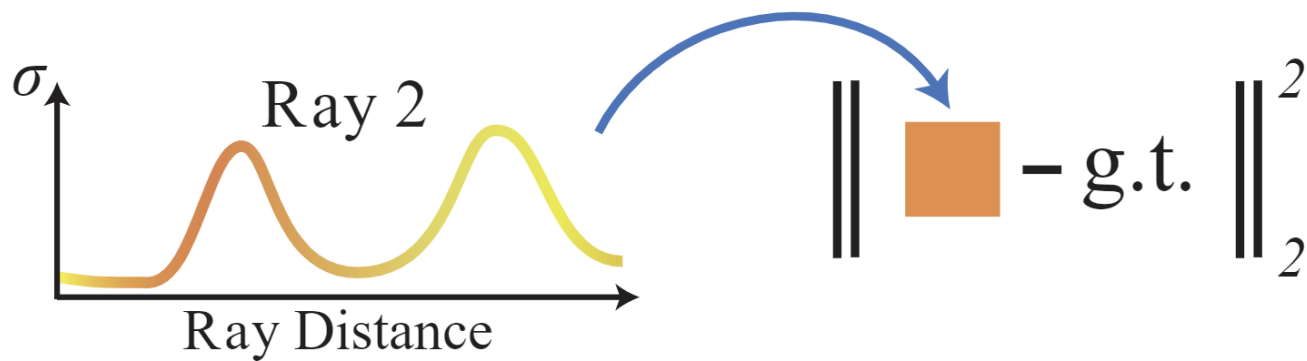
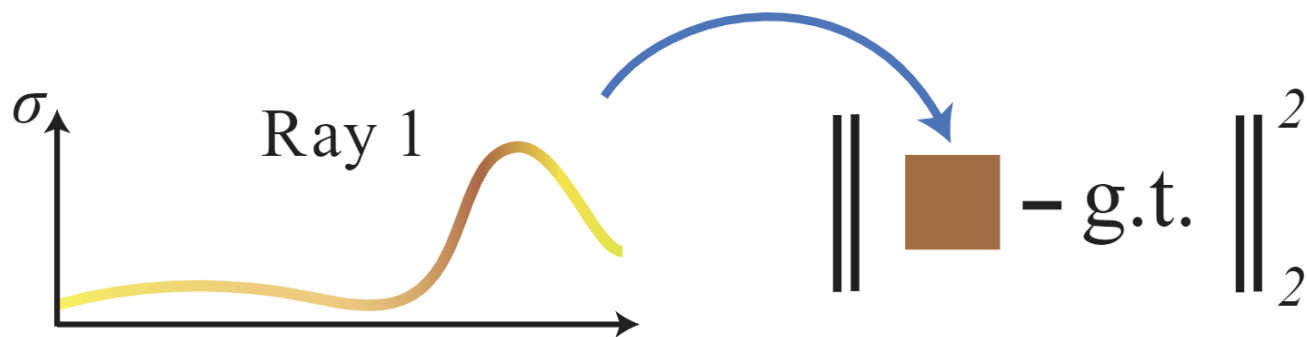
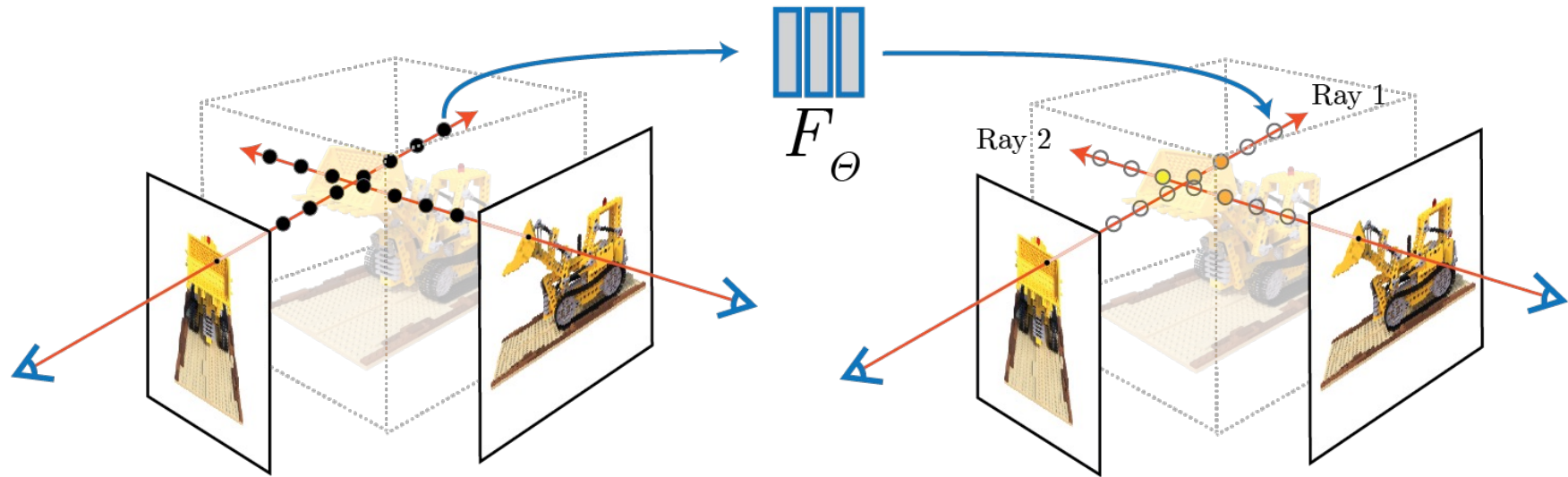


DeepSDF Extensions: NeRF

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



Training NeRFs



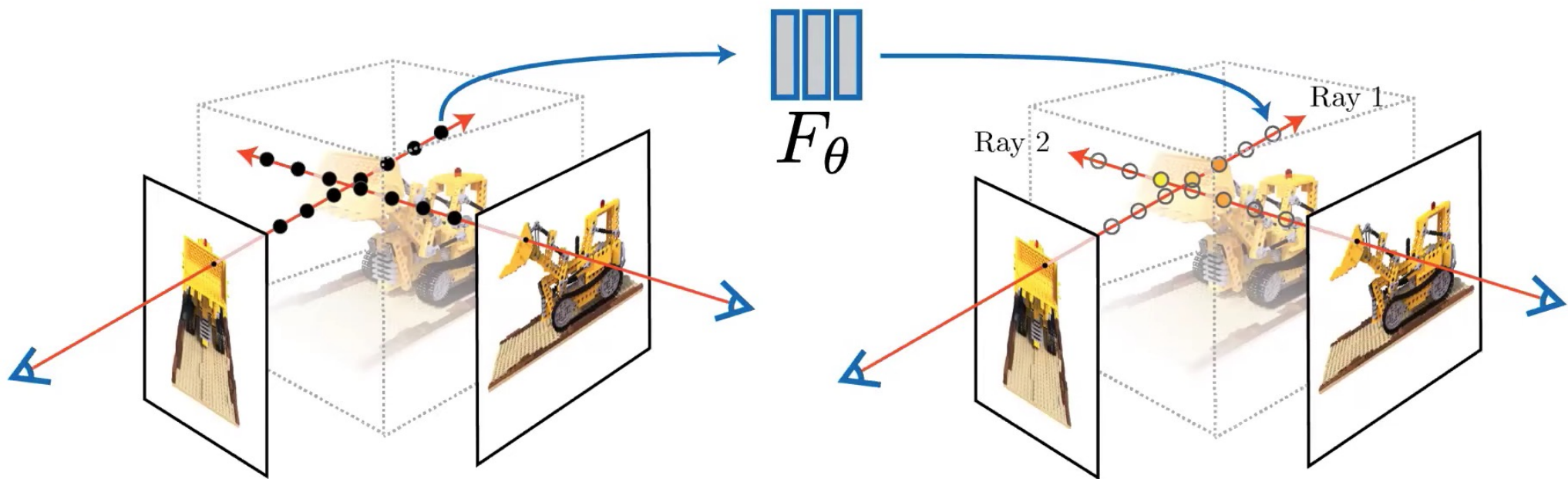


Volume rendering of
MLP colors/densities



Ground truth
image



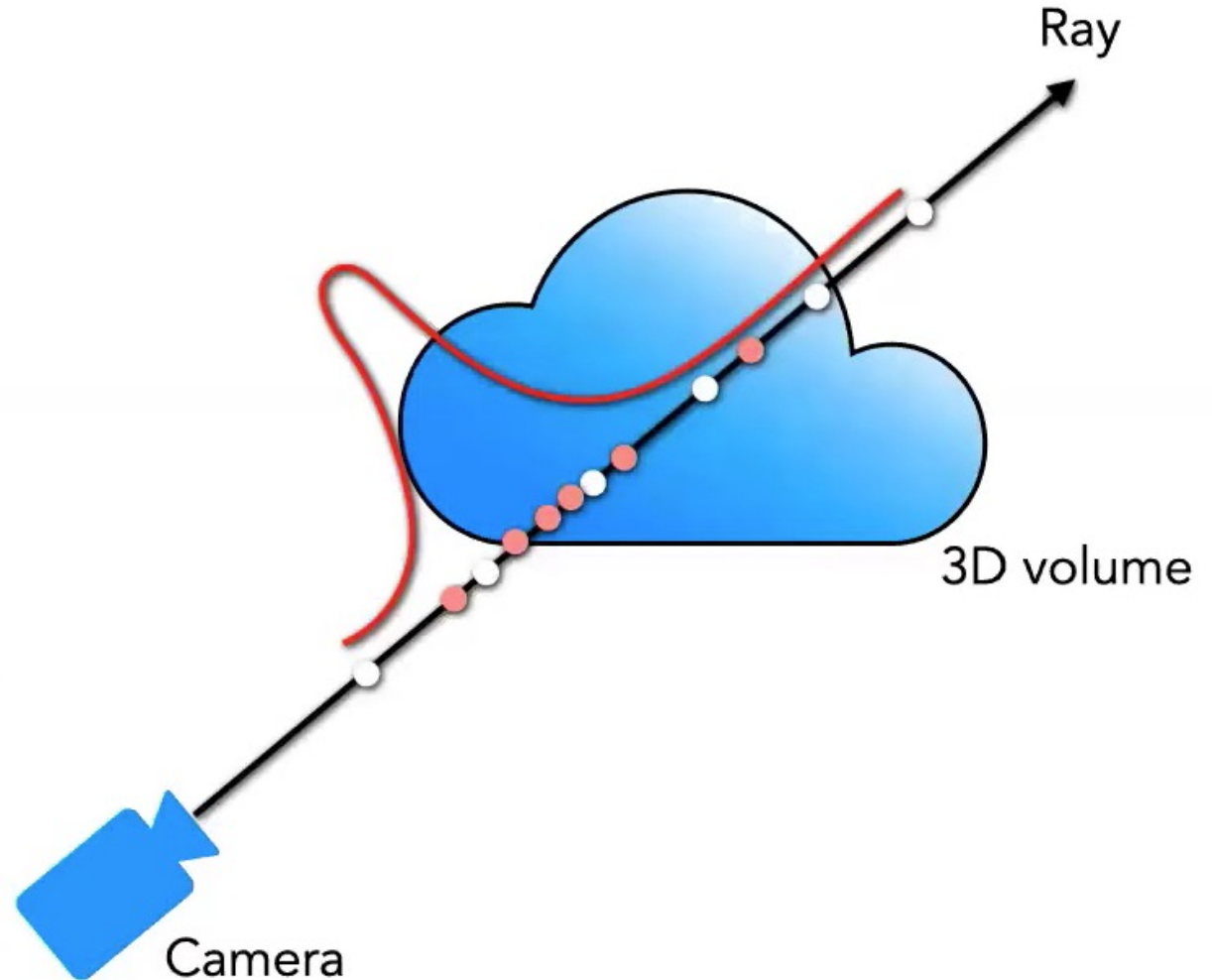


$$\min_{\theta} \sum_i \|\text{render}_i(F_{\theta}) - I_i\|^2$$

Importance Sampling

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

treat weights as probability
distribution for new samples





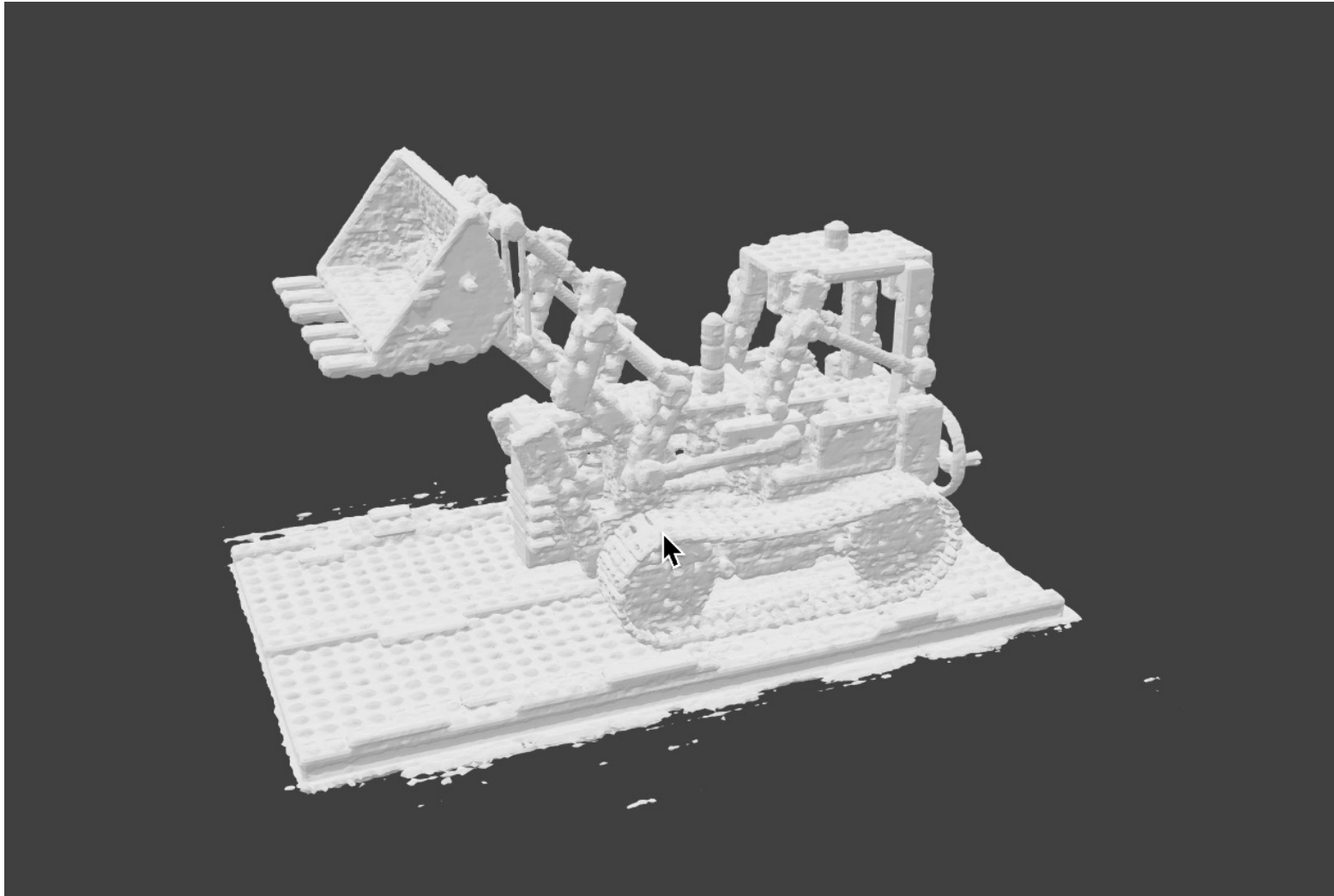


NeRF encodes convincing view-dependent effects using directional dependence



Building 3D models from NeRFs

Apply marching cubes algorithm on NeRF predicted volume density (σ)



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 σ
- 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

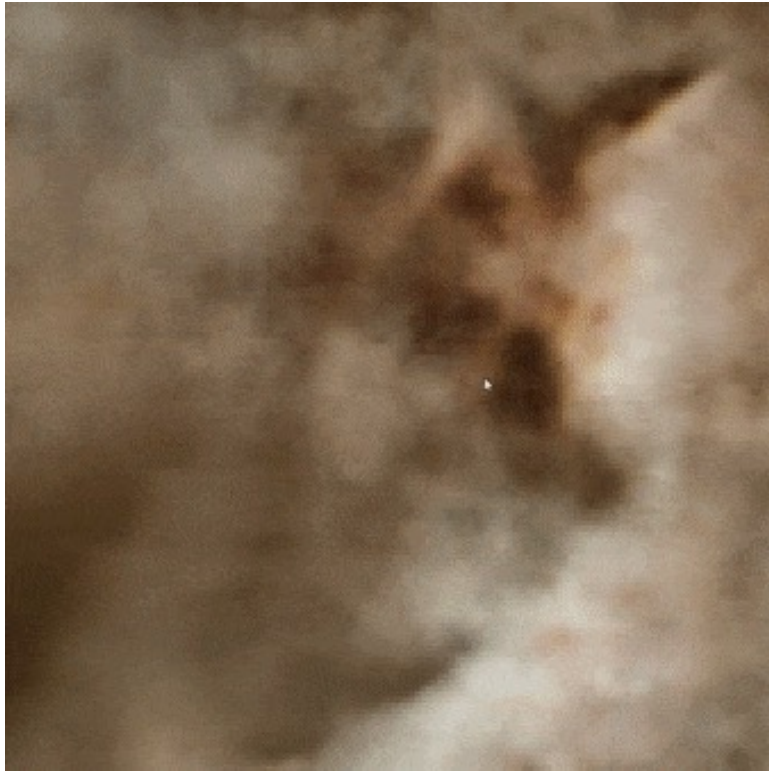
Key limitations of the original NeRF

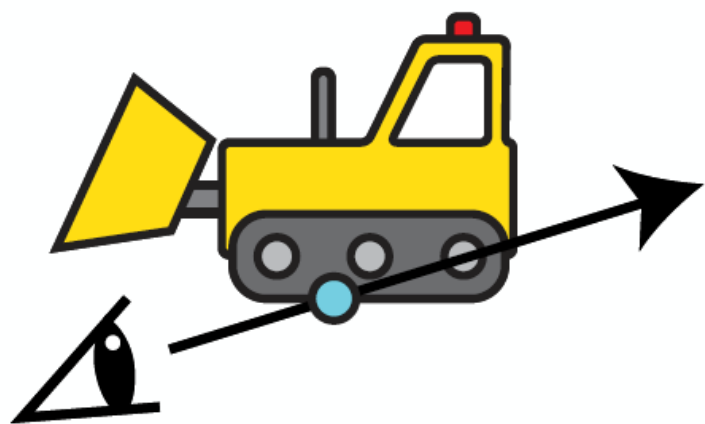
- Very slow in training and inference
- Requires Ground-Truth poses
- Do not generalize to new scenes

Key limitations of the original NeRF

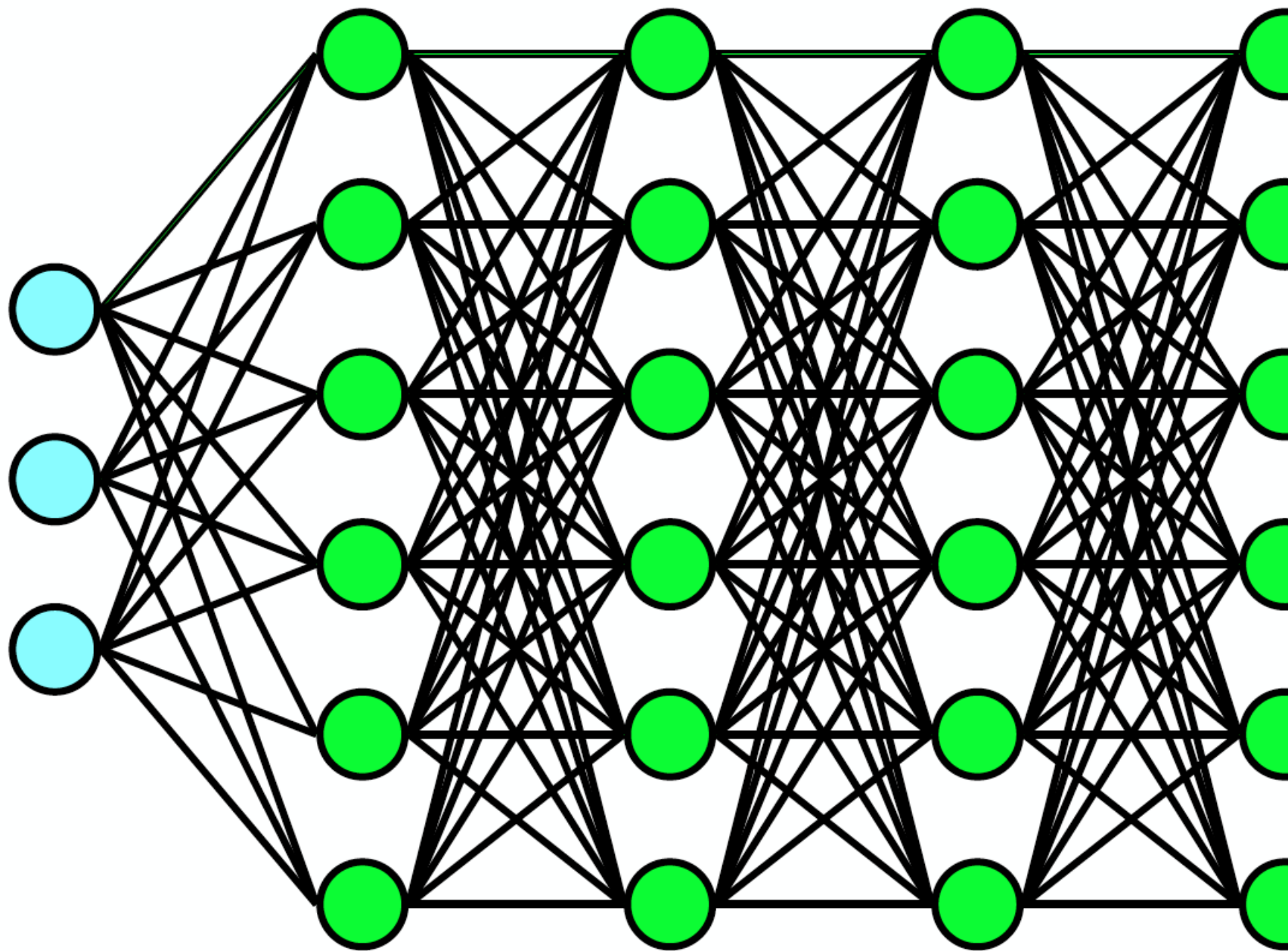
- **Very slow in training and inference**
- Requires Ground-Truth poses
- Do not generalize to new scenes

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





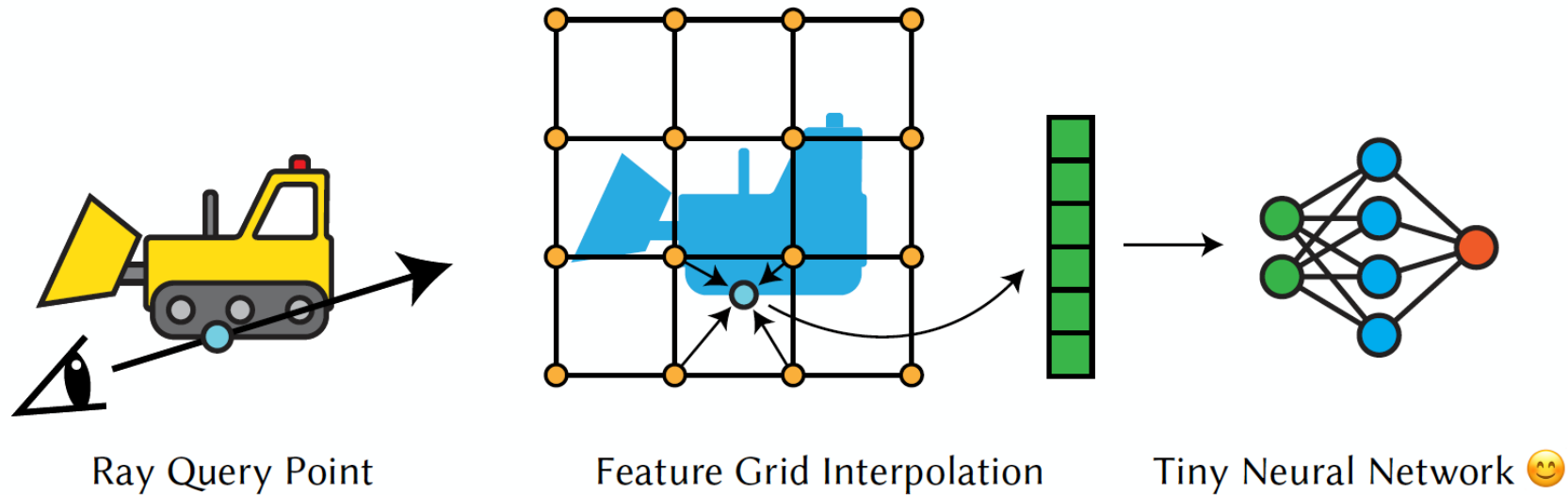
Ray Query Point



Huge Neural Network 🙄

4

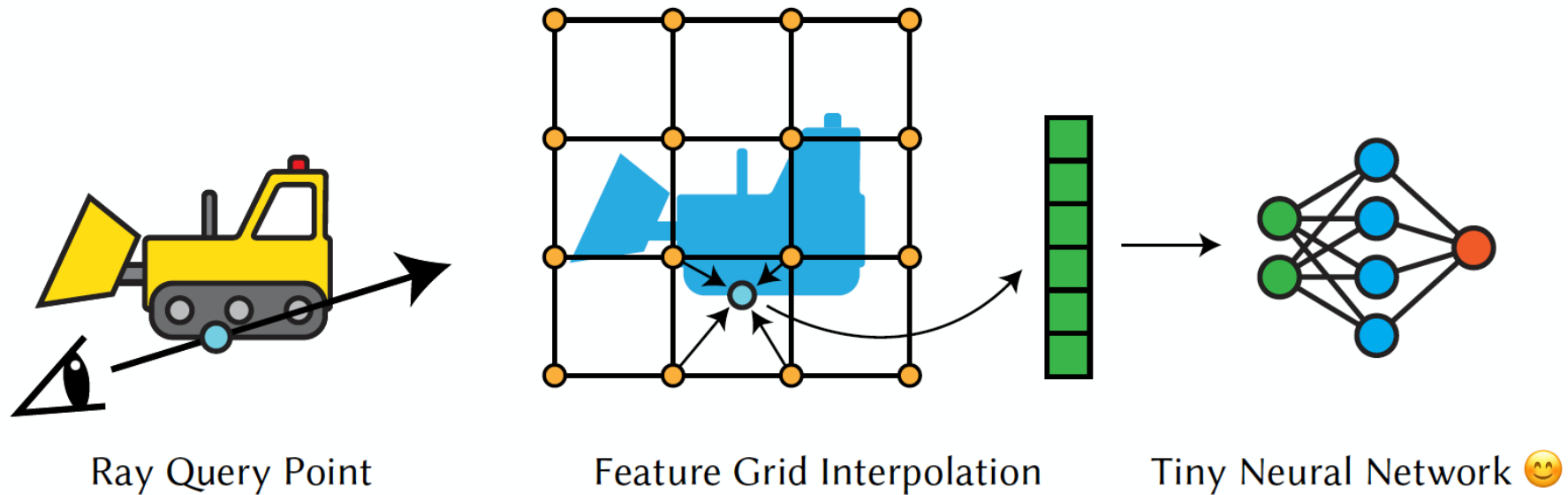
Hybrid representation



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)

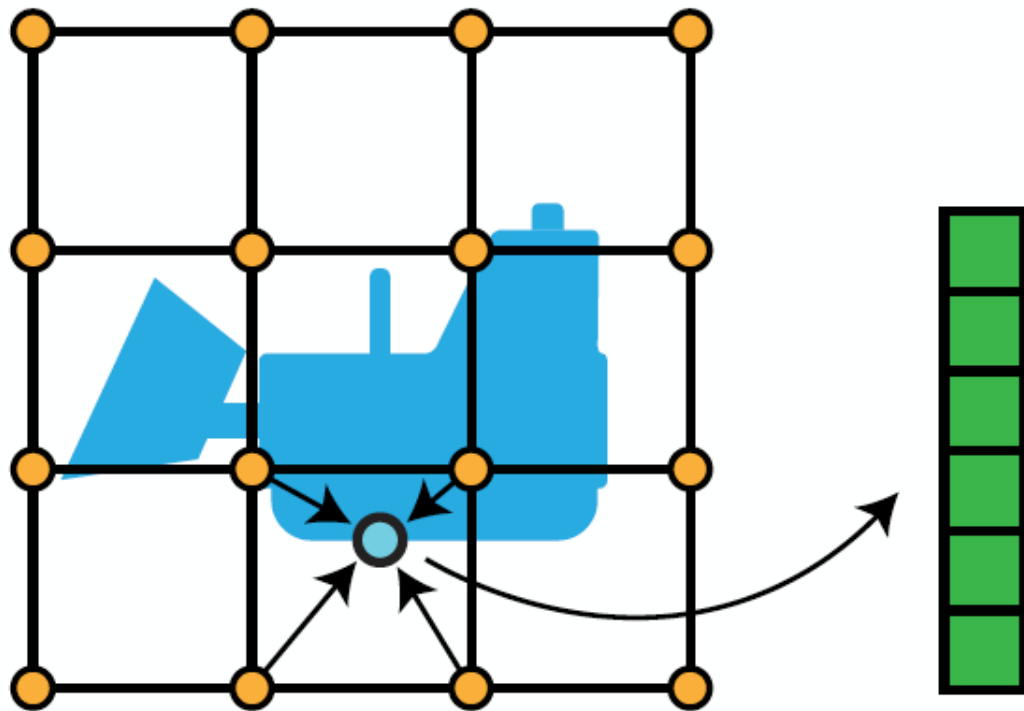
Hybrid representation: It's all about Data Structures!



Why hybrid representation?

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

Uniform Grids



Pros:

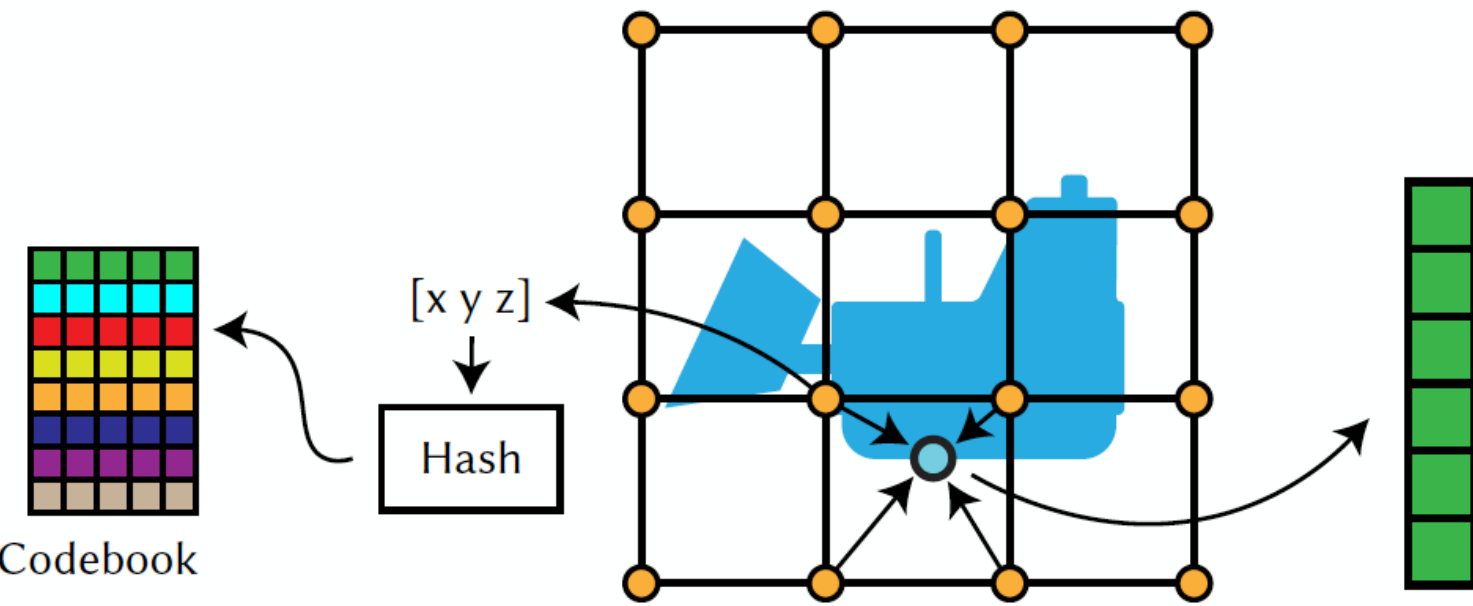
- Easy to implement
- Algorithmically fast access [$O(1)$]
- Established operations like convolutions
- Simple topology

Cons:

- Expensive in memory and bandwidth
- Limited by Nyquist

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

Hash Grids



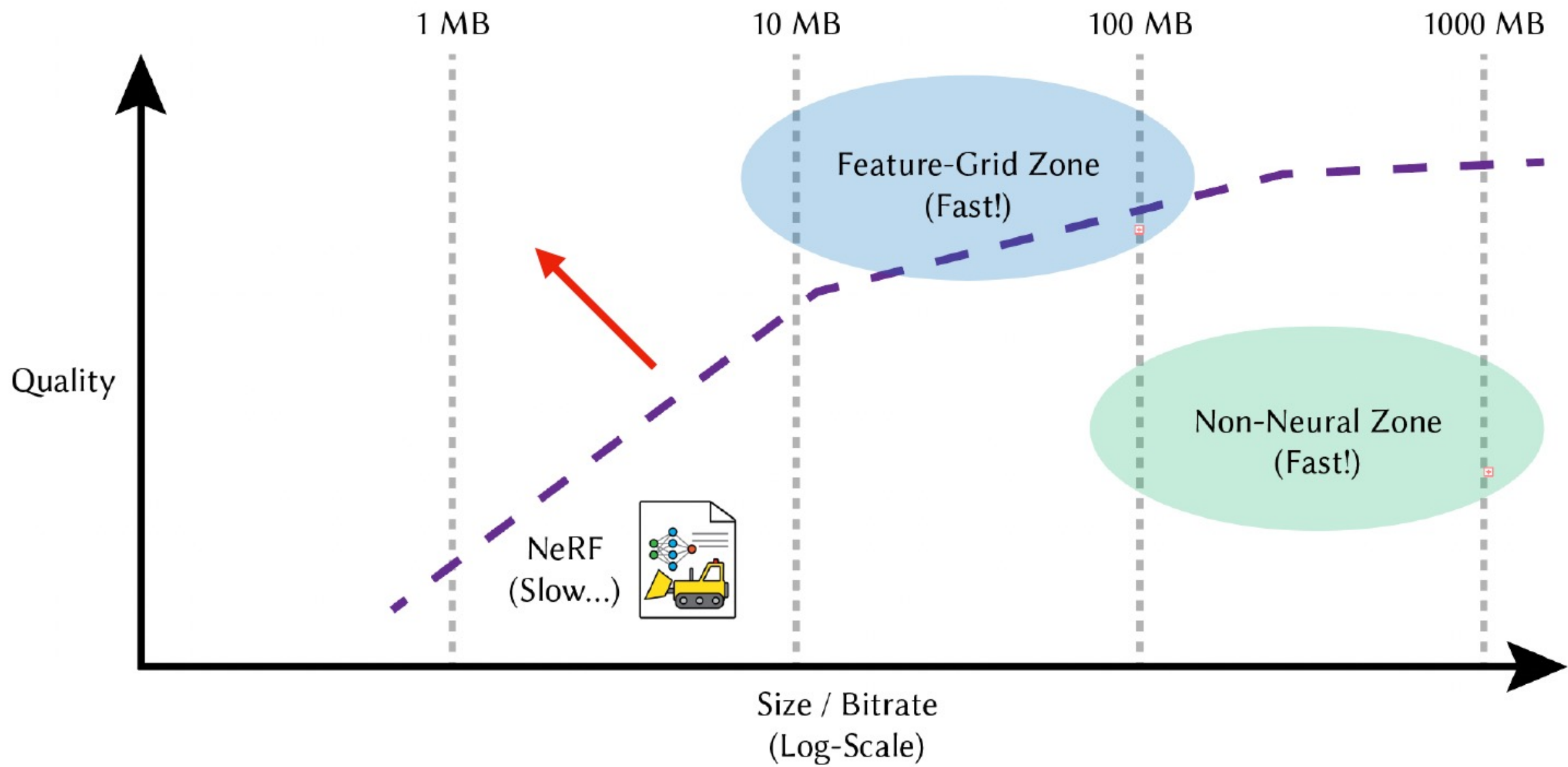
[Instant-NGP (Muller et al.)]

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 🤗

🤗 Carnegie Mellon University

Wei-Chiu Ma 🤢

🤢 Massachusetts Institute of Technology

Antonio Torralba 🤢

Simon Lucey 🤗🤗

🤗 The University of Adelaide

IEEE International Conference on Computer Vision (ICCV), 2021

oral presentation

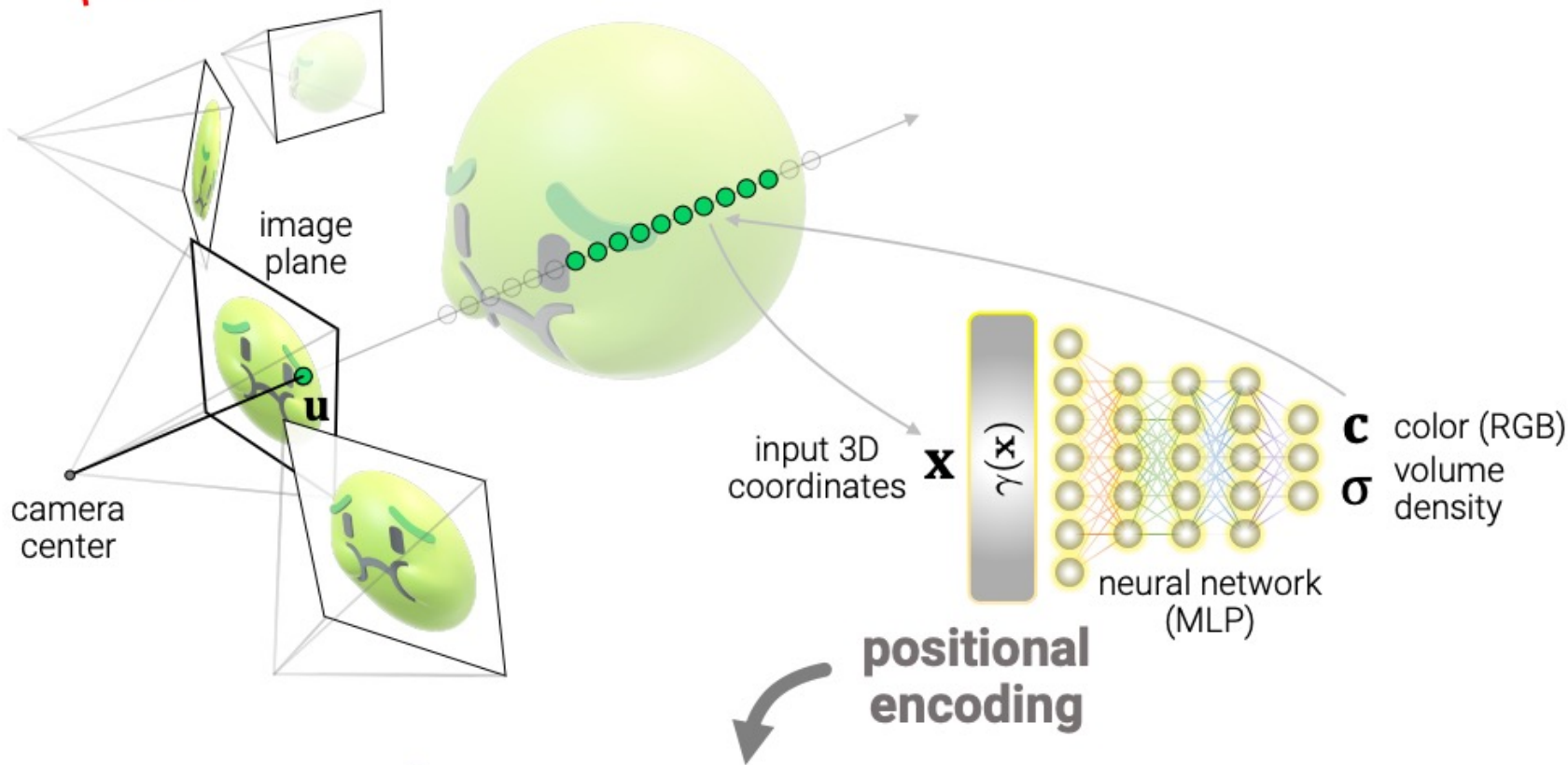
NeRF in a nutshell

We want to optimize the poses as well!

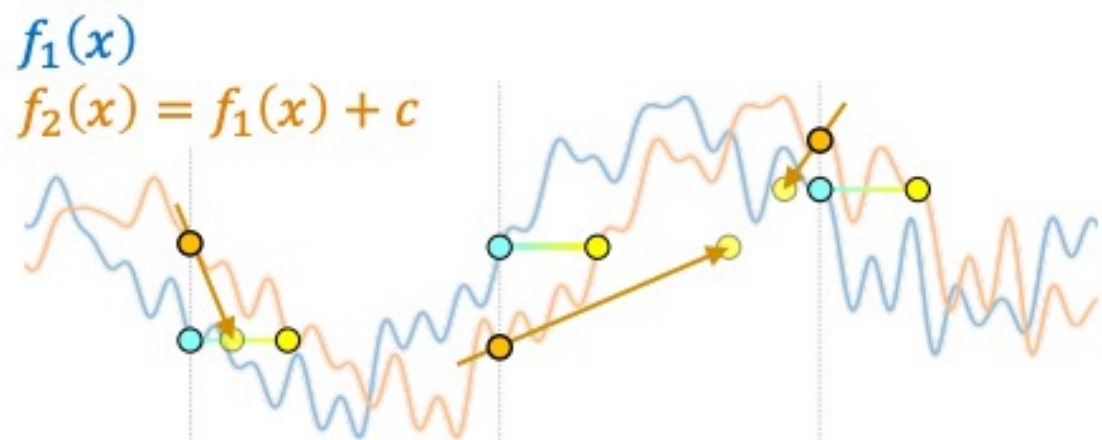
$$\min_{\mathbf{p}_1, \dots, \mathbf{p}_M, \Theta}$$

$$\sum_{i=1}^{\text{frames } M} \sum_{\mathbf{u}} \left\| \hat{\mathcal{I}}(\mathbf{u}; \mathbf{p}_i, \Theta) - \mathcal{I}_i(\mathbf{u}) \right\|_2^2$$

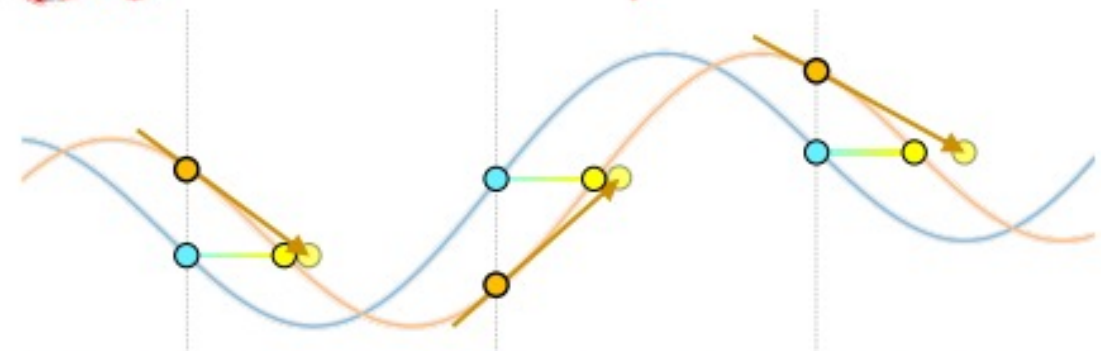
RGB (rendered) camera pose network params. RGB



- ✓ encourages representation with high frequency
- ✗ detrimental to gradient-based **registration!!**



⊗ gets stuck in suboptimal solutions



✓ smooth signals → coherent updates

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

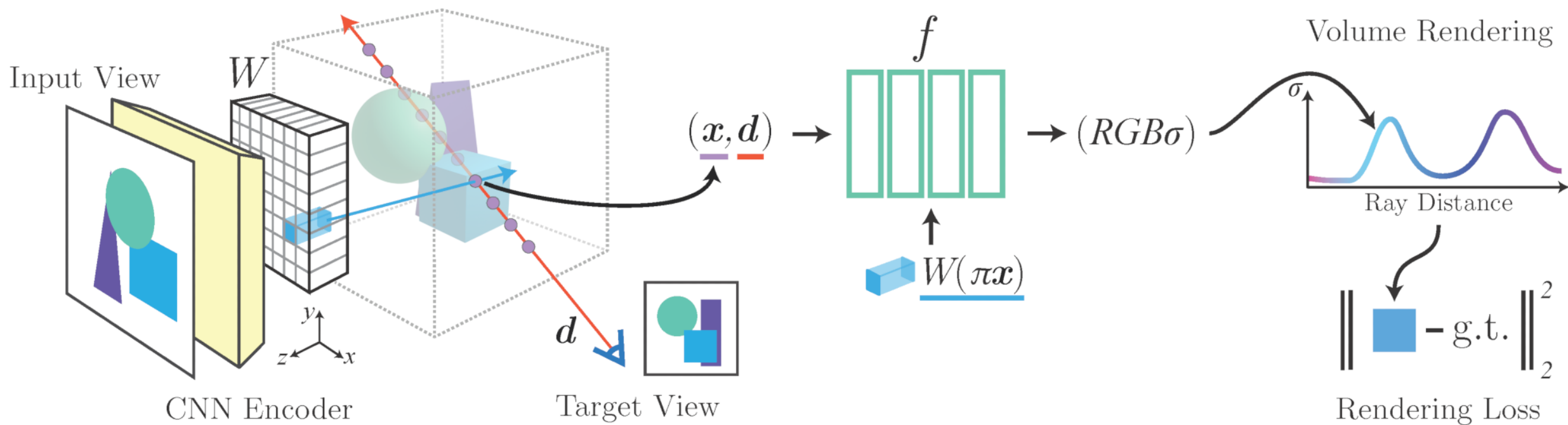
pixelNeRF

Neural Radiance Fields from One or Few Images

CVPR 2021

Alex Yu Vickie Ye Matthew Tancik Angjoo Kanazawa

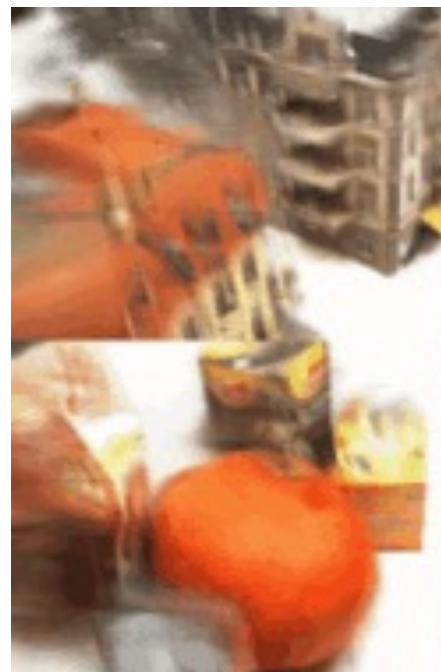
UC Berkeley



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.