# 3D-aware Synthesis

## Jun-Yan Zhu

16-726, Spring 2023
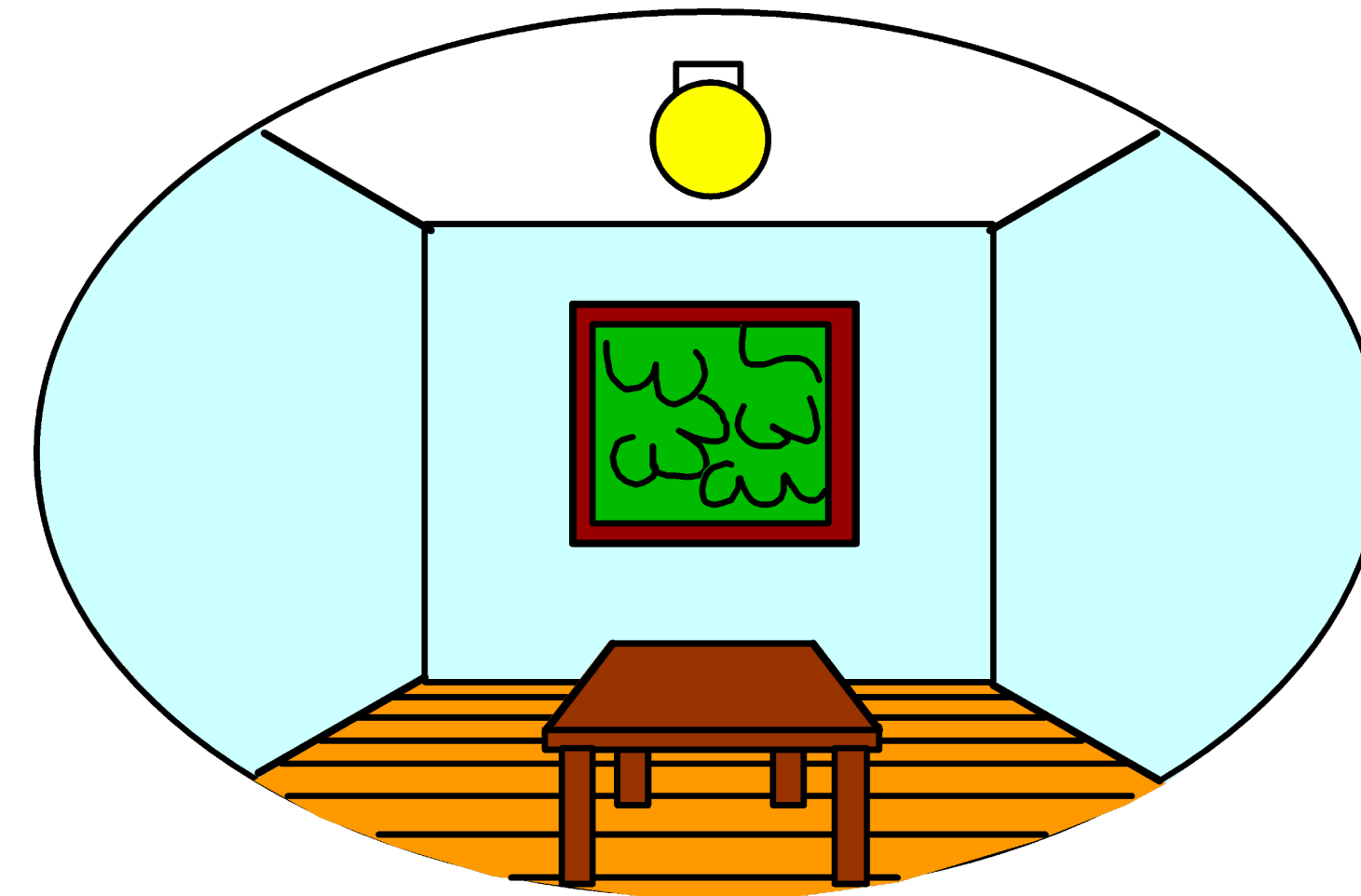
# What do we see?
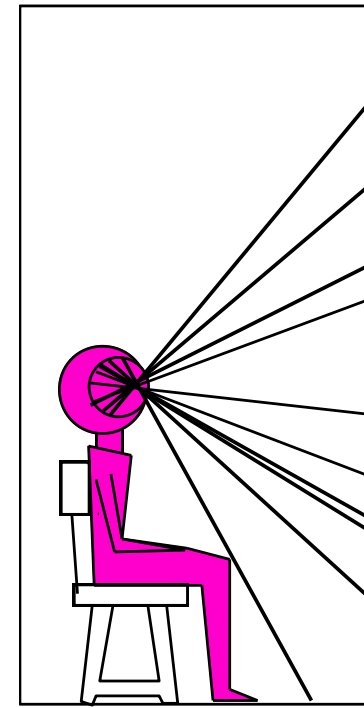
**3D world**
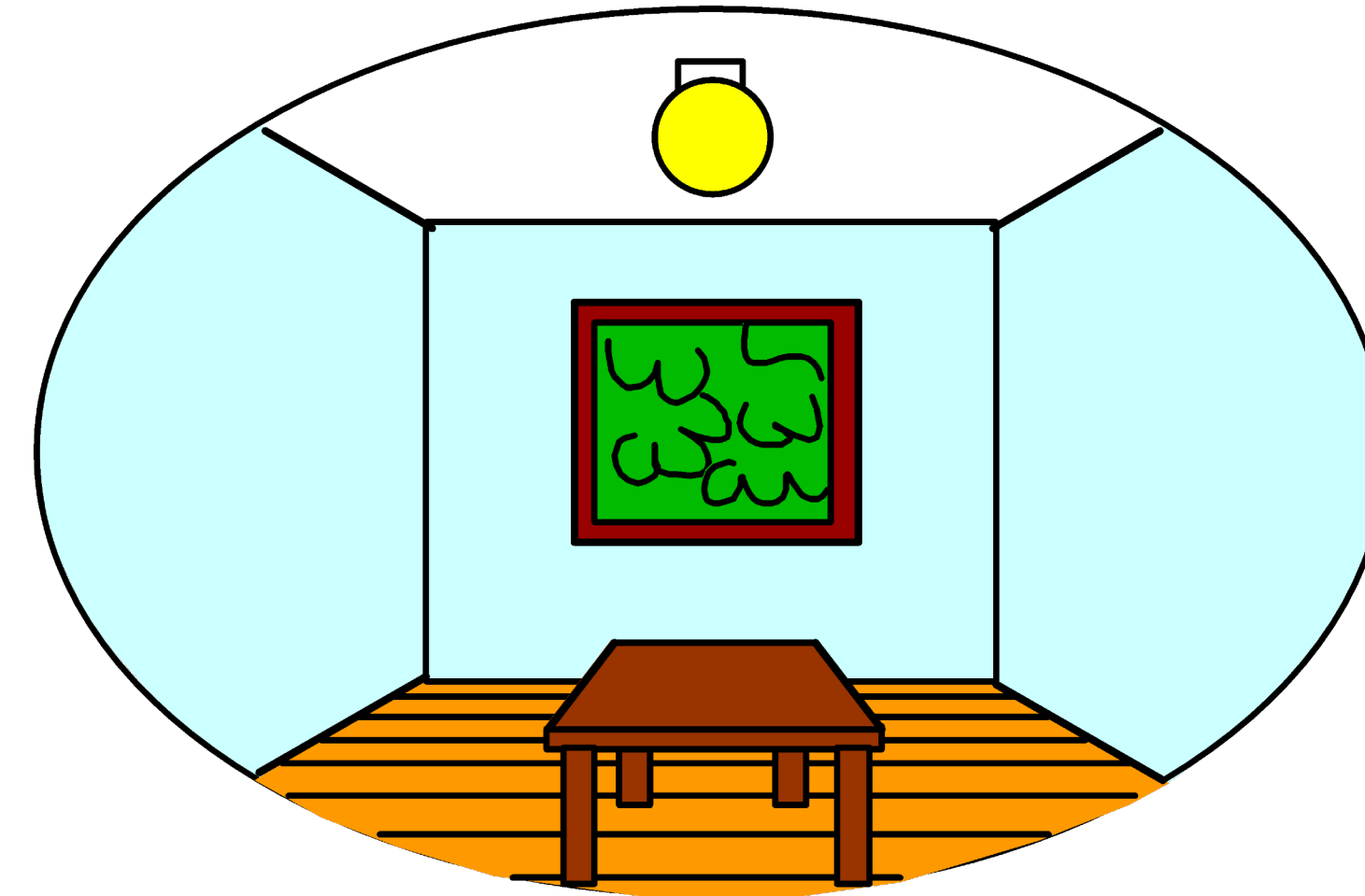
**2D image**

Point of observation

# What do we see?

**3D world**

**2D image**

Painted
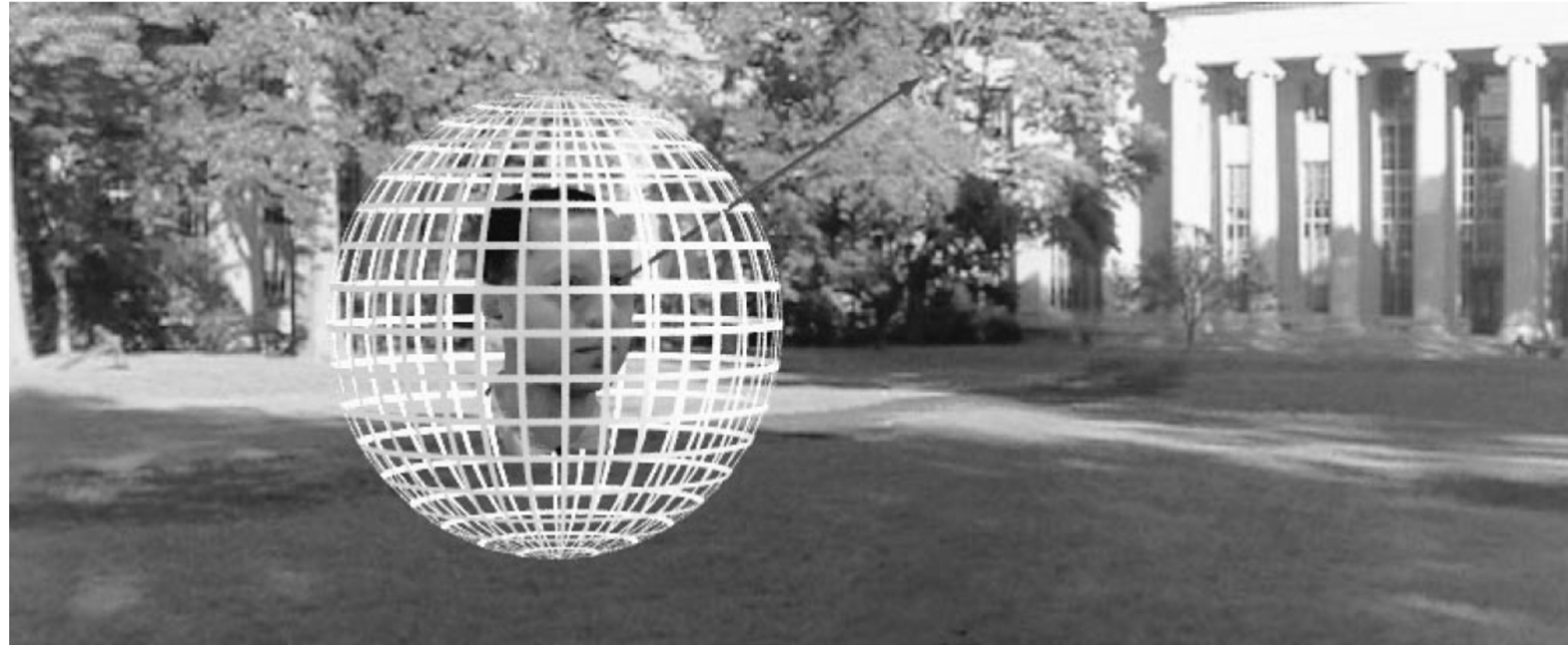backdrop

# The Plenoptic Function



**Figure by Leonard McMillan**

- Q: What is the set of all things that we can ever see?

- A: The Plenoptic Function (Adelson & Bergen)

- Let's start with a stationary person and try to parameterize <u>everything</u> that she or he can see…

# Grayscale snapshot



$$P(\theta,\phi)$$

- is intensity of light

  - Seen from a single view point

  - At a single time

  - Averaged over the wavelengths of the visible spectrum

- (can also do *P(x,y),* but spherical coordinate are nicer)

# Color snapshot



$$P(\theta,\phi,\lambda)$$

- is intensity of light

  - Seen from a single view point
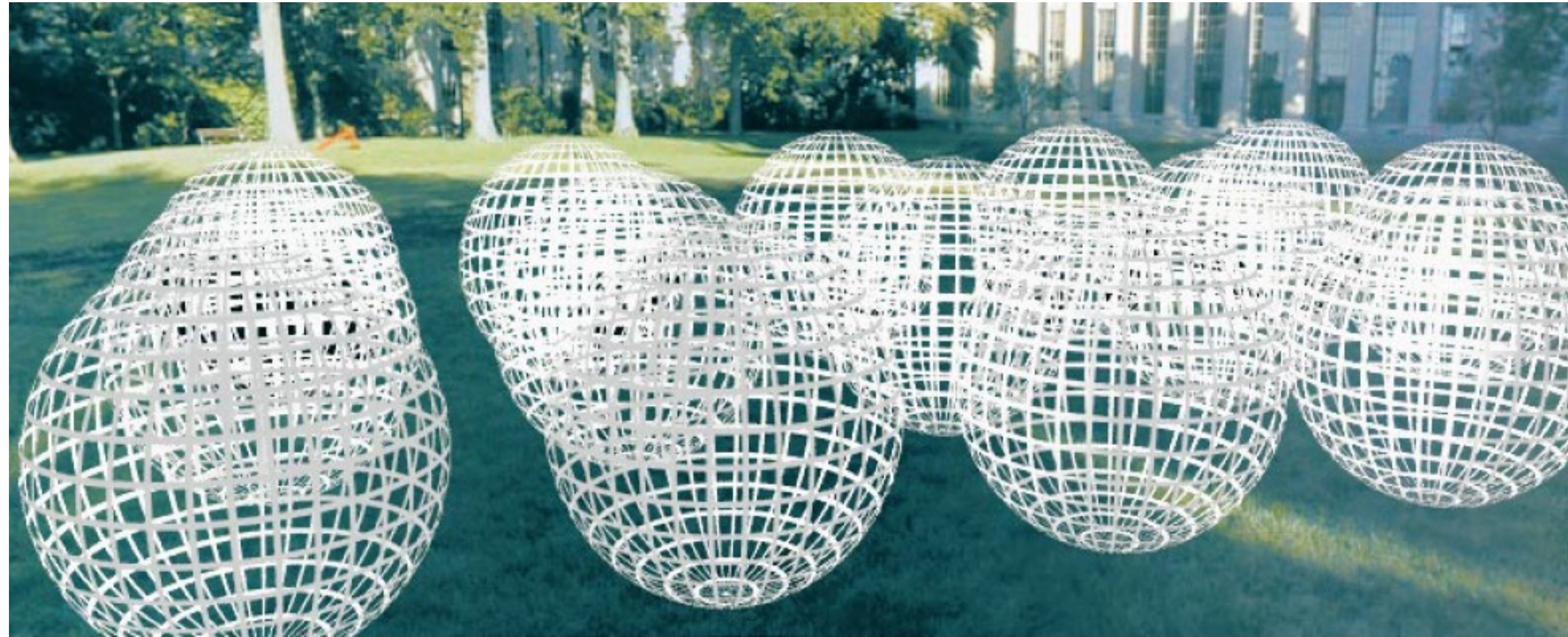
  - At a single time

  - As a function of wavelength

# A movie



$$P(\theta,\phi,\lambda,t)$$

- is intensity of light

  - Seen from a single view point

  - Over time

  - As a function of wavelength
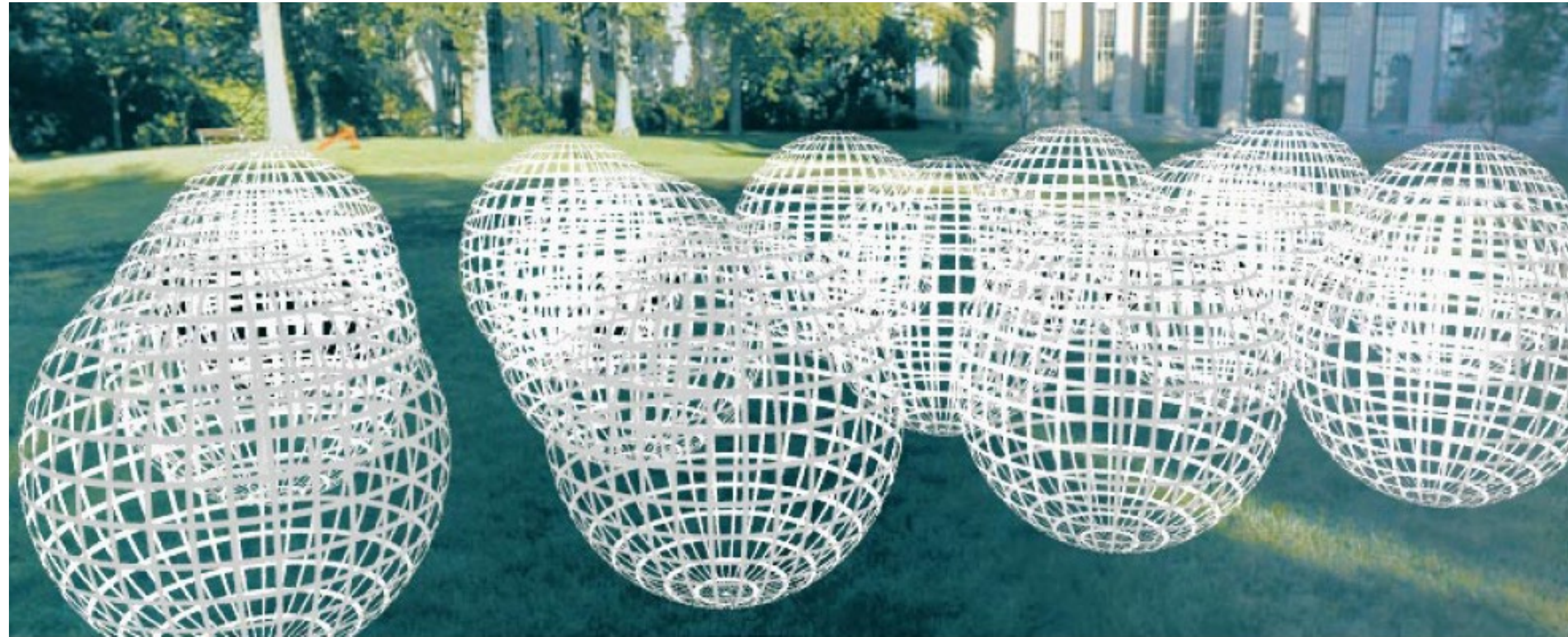
# Holographic movie



$$P(\theta, \phi, \lambda, t, V_X, V_Y, V_Z)$$

- is intensity of light

  - Seen from ANY viewpoint

  - Over time
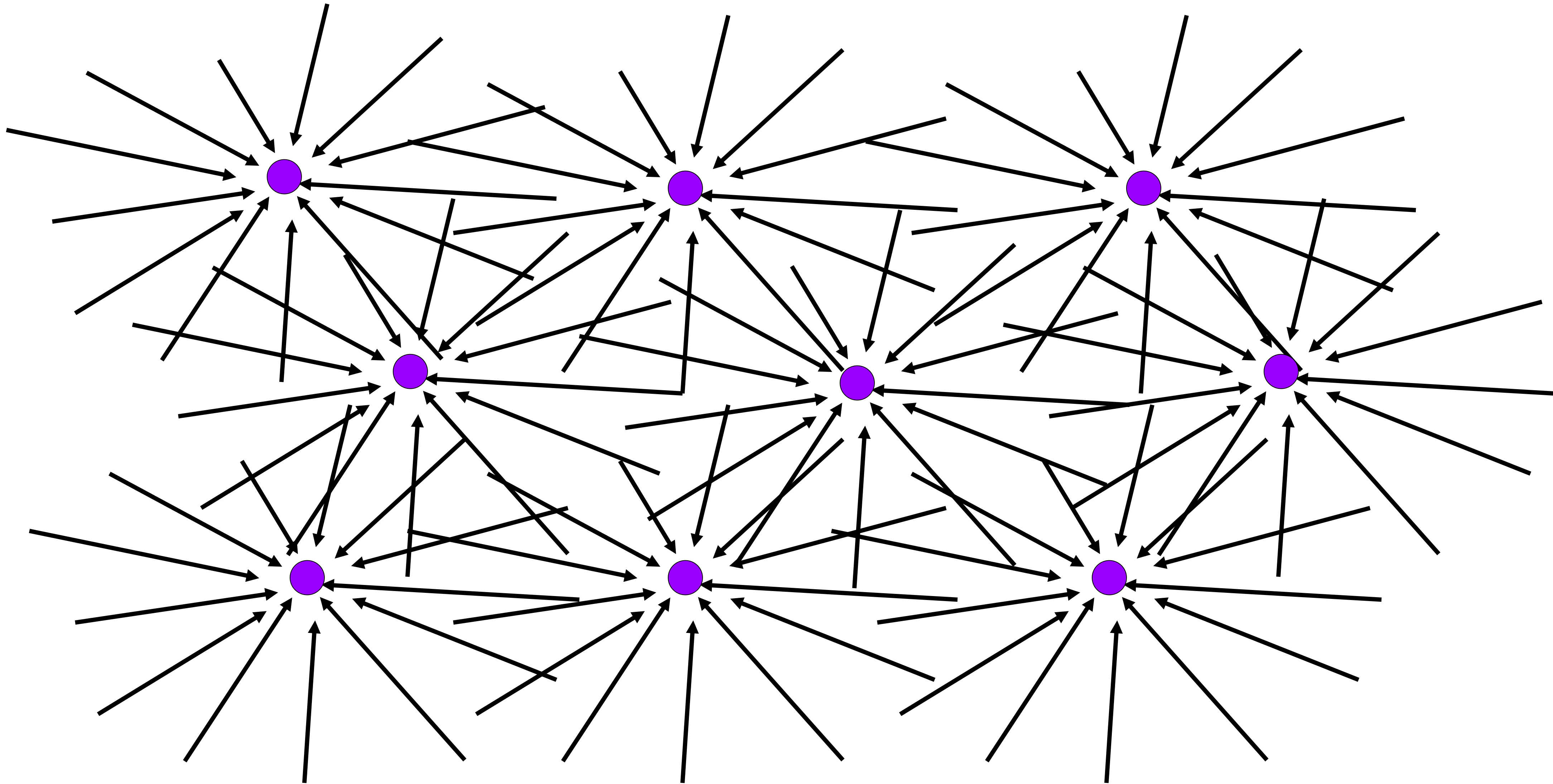
  - As a function of wavelength

# The Plenoptic Function



$$P(\theta, \phi, \lambda, t, V_X, V_Y, V_Z)$$

- Can reconstruct every possible view, at every moment, from every position, at every wavelength

- Contains every photograph, every movie, everything that anyone has ever seen! it completely captures our visual reality!   Not bad for a function…

# Sampling Plenoptic Function (top view)



**Just lookup -- Quicktime VR**

# QuickTime VR

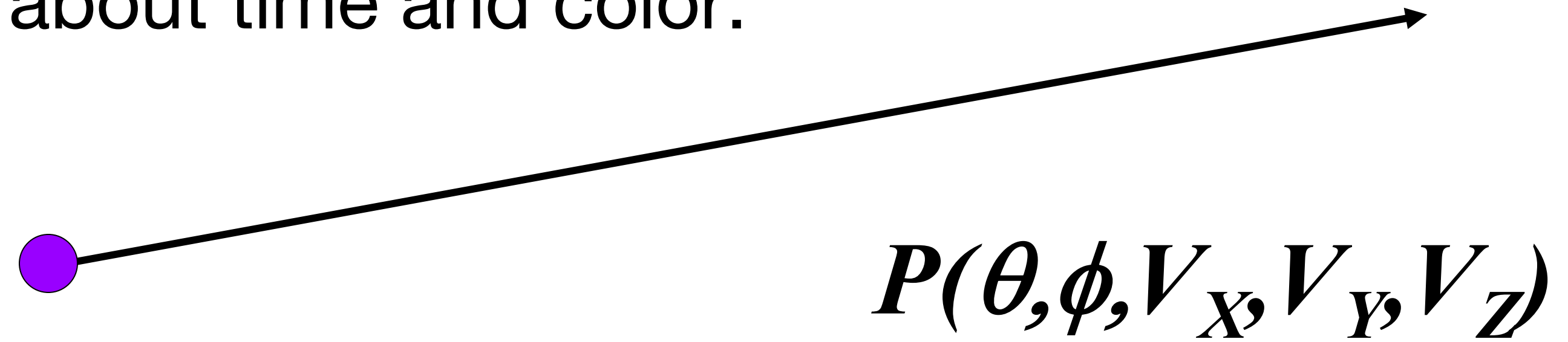**Panoramic image**



**Perspective Warp**

# QuickTime VR



Quicktime VR: An image-based approach to virtual environment navigation. Shenchang Eric Chen. SIGGRAPH 1995
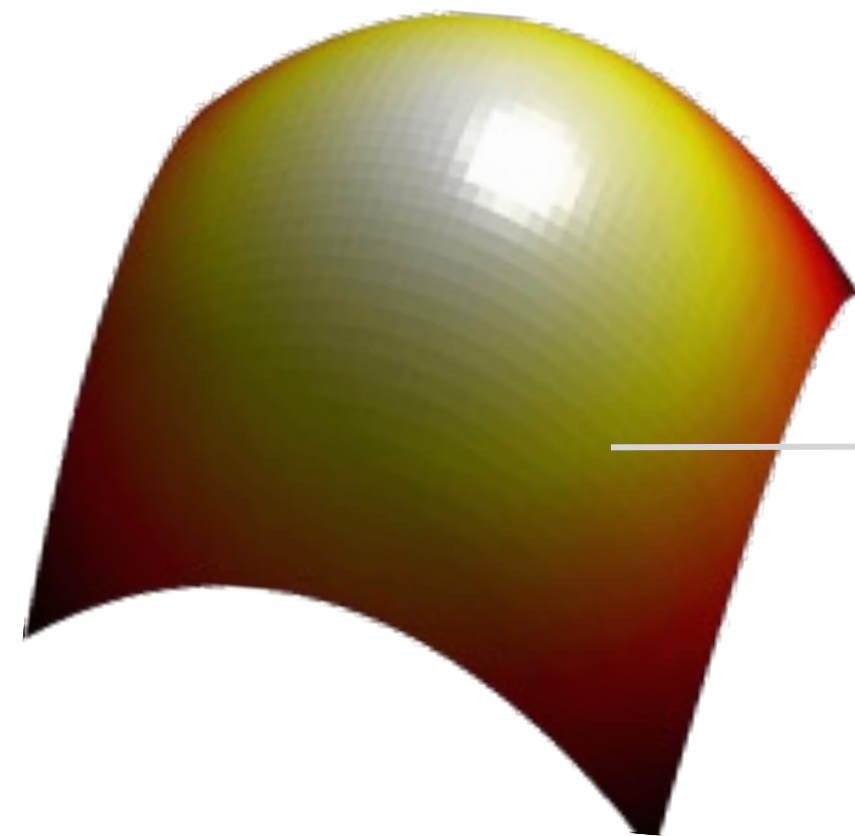
# Ray

- Let's not worry about time and color:

$$P(\theta, \phi, V_X, V_Y, V_Z)$$

- 5D

  - 3D position

  - 2D direction

# How can we use this?
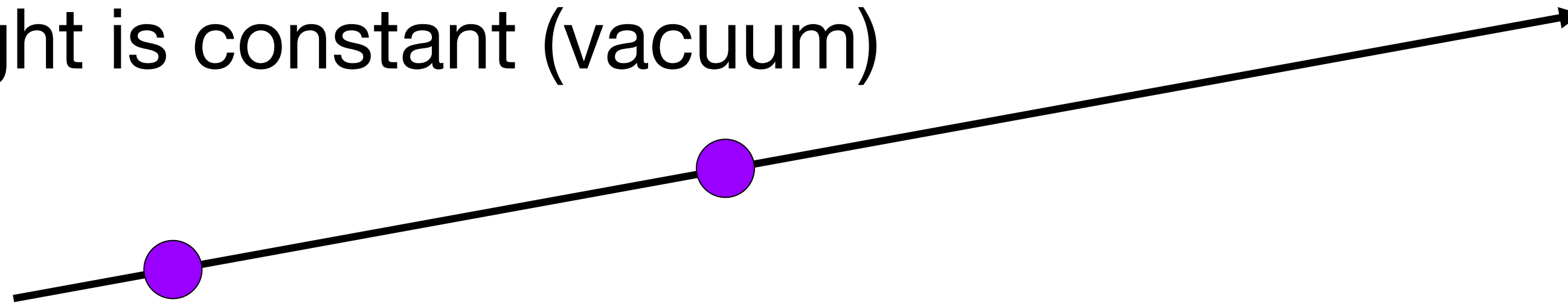
Lighting

Surface

No Change in

Radiance

Camera

# Ray Reuse

- Infinite line
- Assume light is constant (vacuum)

- 4D
  - 2D direction
  - 2D position
  - non-dispersive medium
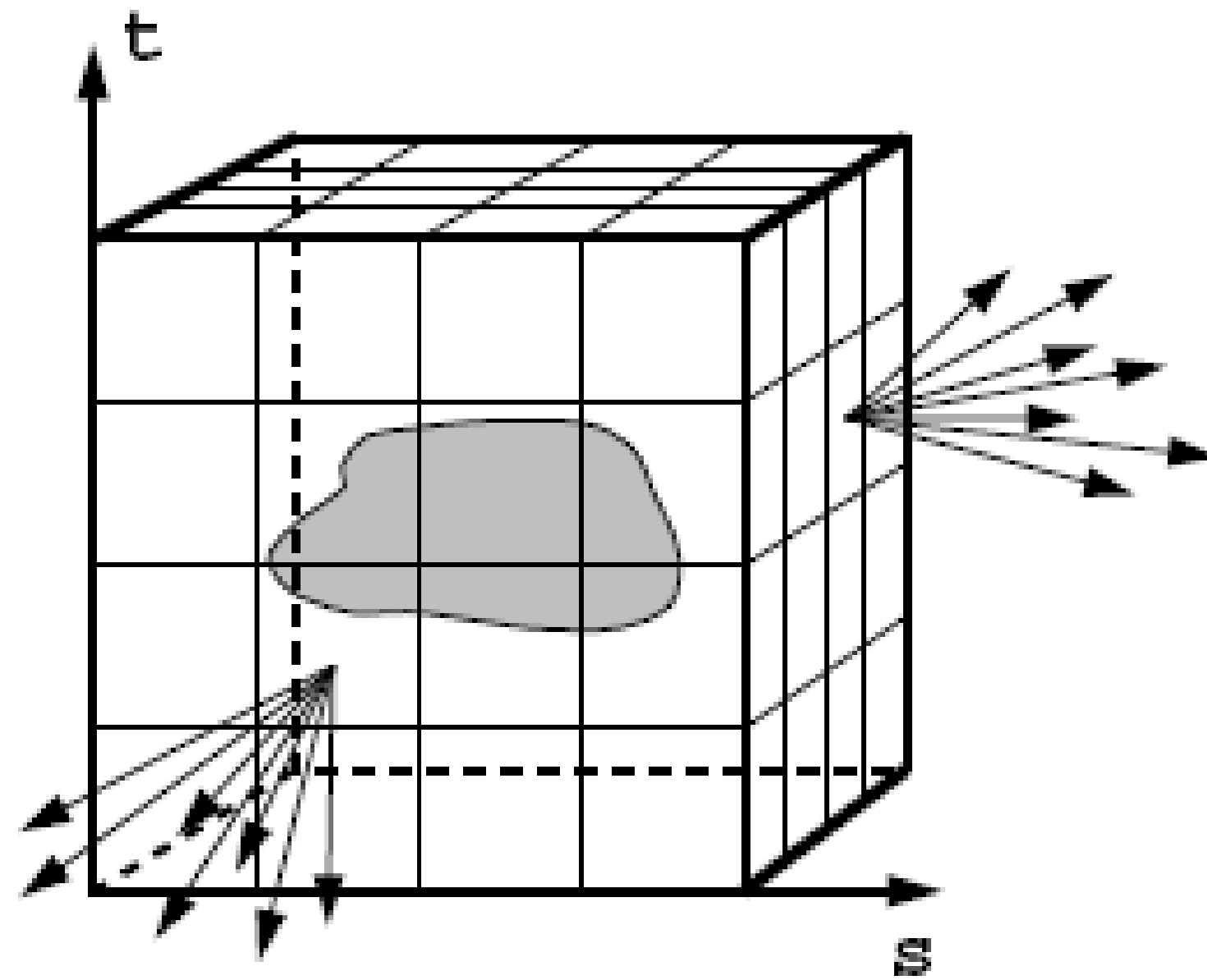
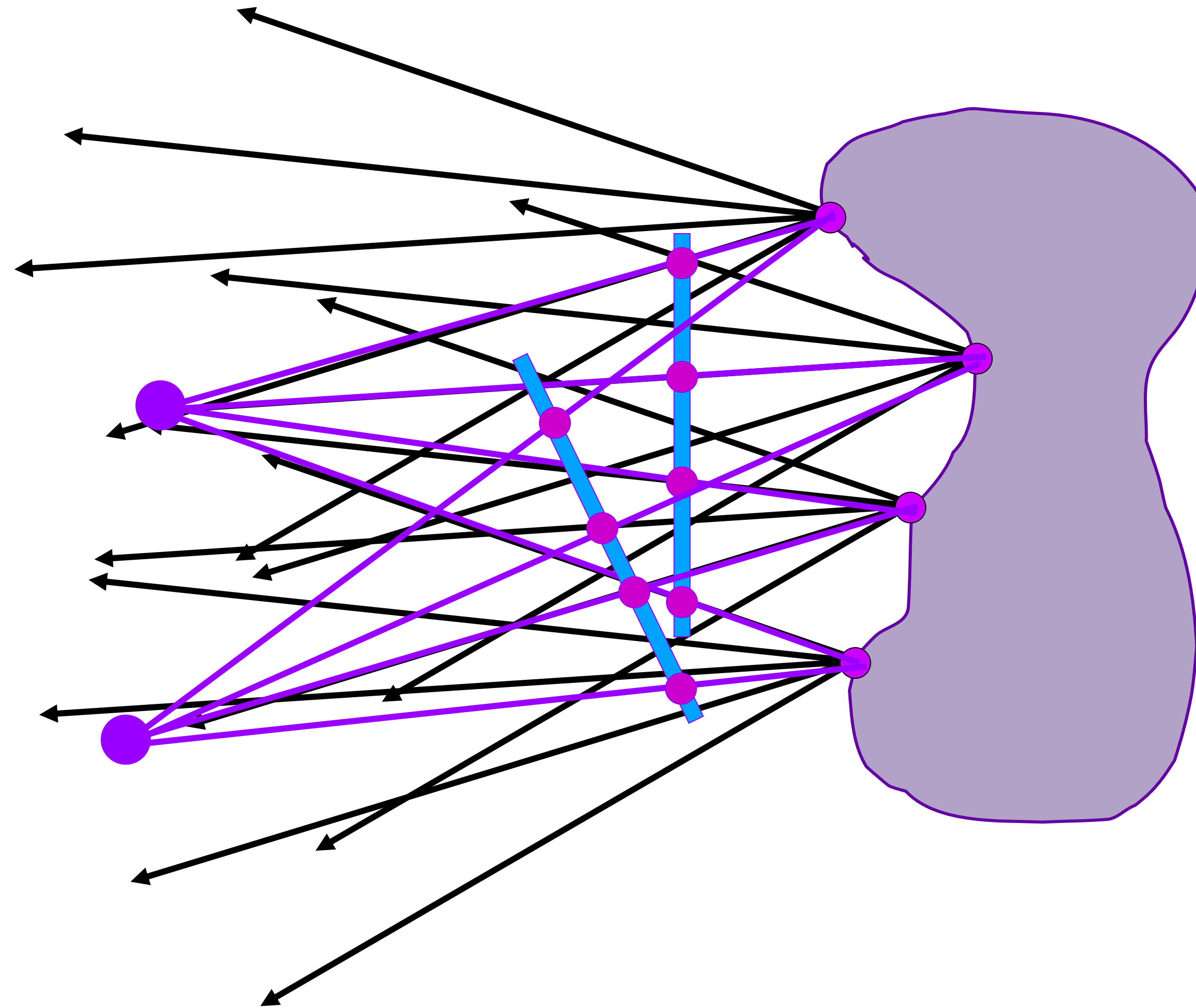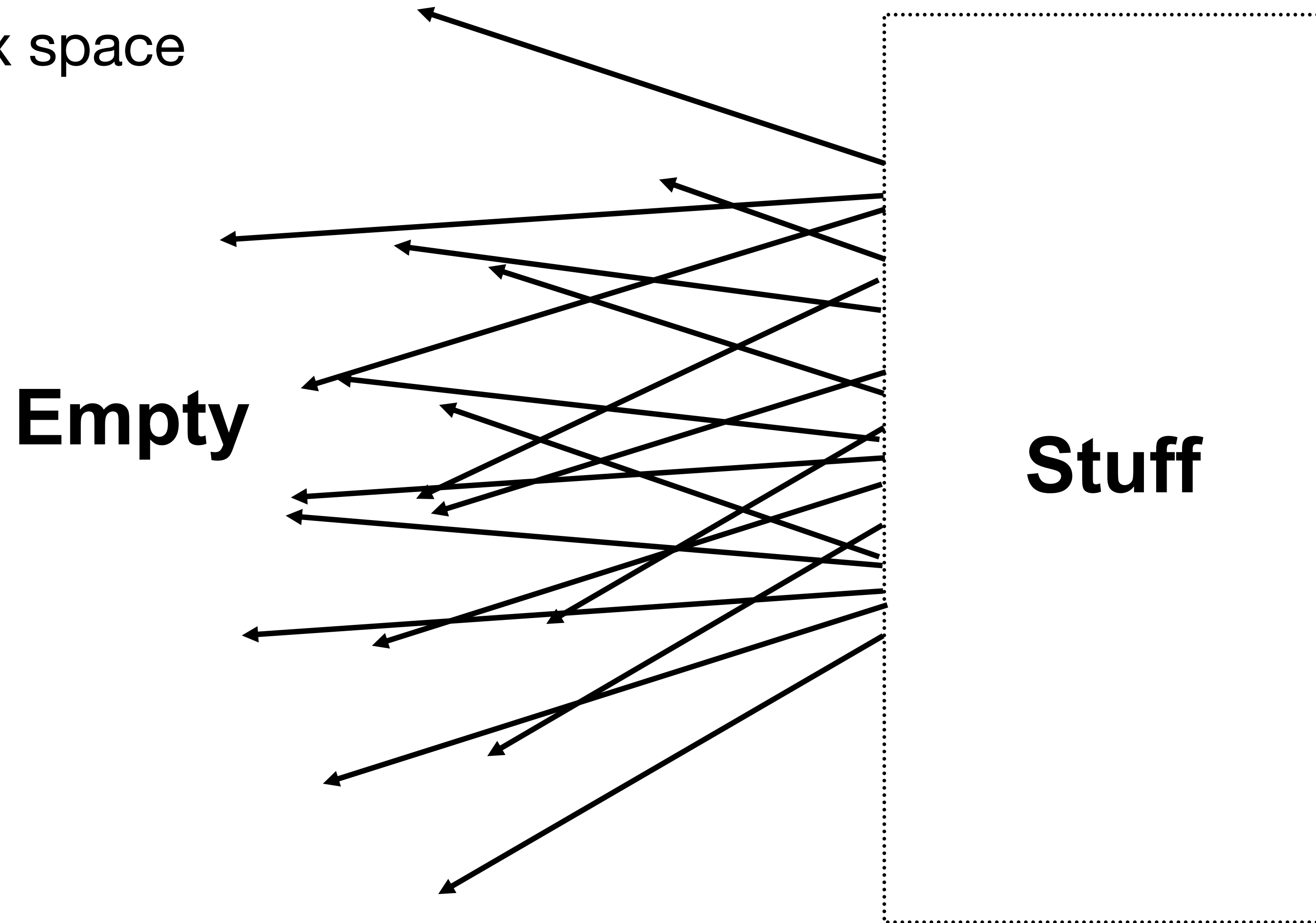# Only need plenoptic surface



Figure 1: The surface of a cube holds all the radiance information due to the enclosed object.

# Synthesizing novel views

# Lumigraph / Lightfield

- Outside convex space

**Empty**

**Stuff**

- 4D

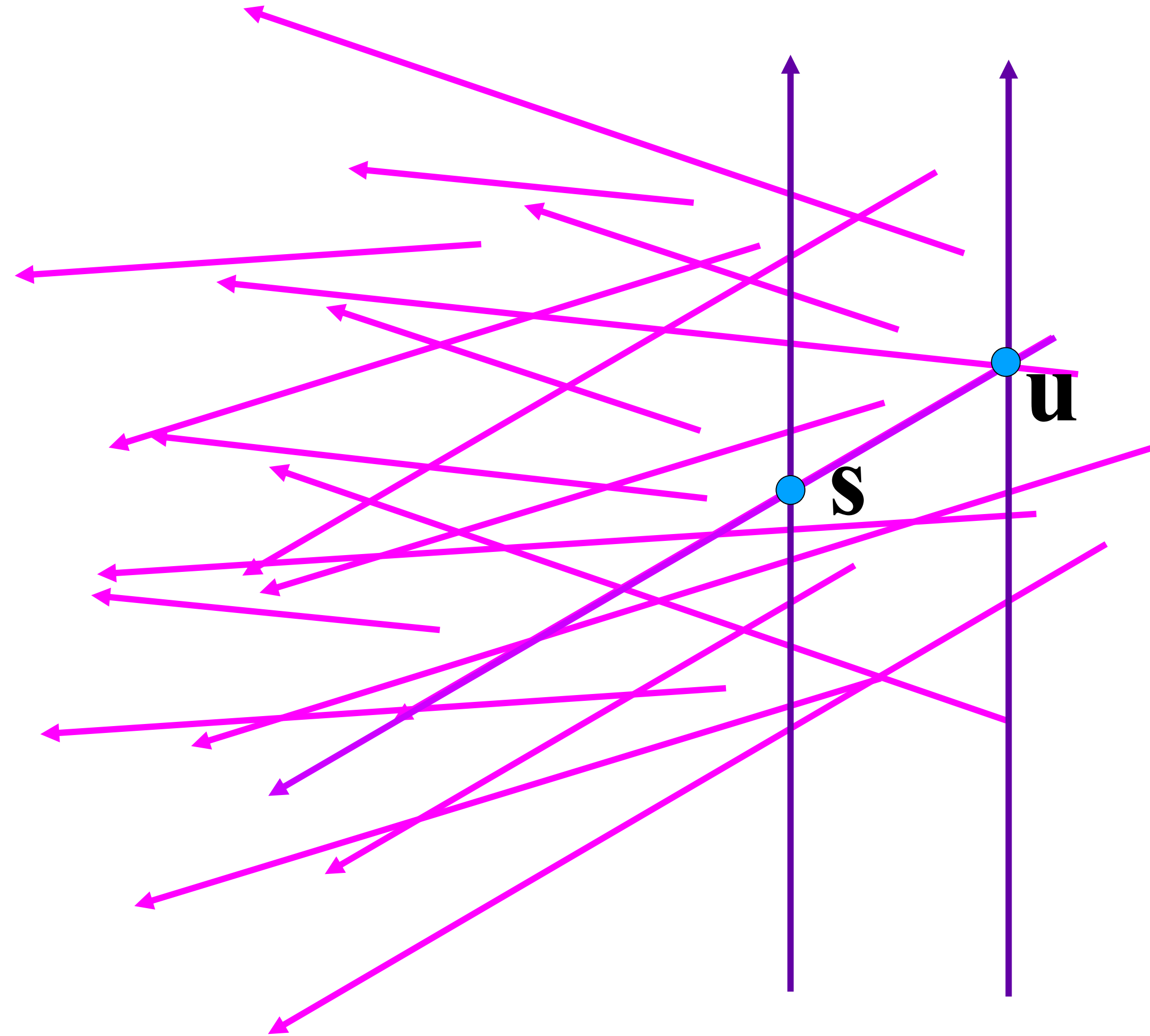# Lumigraph - Organization

- 2D position

- 2D direction



$\theta$

**s**

# Lumigraph - Organization

2D position
2D position

2 plane parameterization

s

u

# Lumigraph - Organization

2D position

2D position

2D position

2 plane parameterization

# Lumigraph - Organization

**Hold s,t constant**

**Let u,v vary**

**An image**



s,t     u,v

# Lumigraph / Lightfield

s

t

u

v

u

v

s

t

# Capture Light Field

## Idea 1

- Move camera carefully over s, t plane
- Grantry
  - see Lightfield paper
  
  [Marc Levoy and Pat Hanrahan]



s,t    u,v

# Capture Light Field

## Idea 2

- Move camera anywhere
- Interpolation over irregular samples
  –see Lumigraph paper
  [Gortler, Grzeszczuk, Szeliski,
  Cohen]

**s,t**      **u,v**

# Novel View Synthesis

For each output pixel
- determine s,t,u,v
- use closest discrete RGB

OR
- interpolate near values

**s**        **u**

# Interpolation

- Nearest neighbor
  - closest s
  - closest u
  - draw it

- Blend 16 nearest
  - quadrilinear interpolation

**s**

**u**

# Stanford multi-camera array



- 640 × 480 pixels × 30 fps × 128 cameras

- synchronized timing

- continuous streaming

- flexible arrangement

# Light field photography using a handheld plenoptic camera

*Ren Ng, Marc Levoy, Mathieu Brédif,*
*Gene Duval, Mark Horowitz and Pat Hanrahan*

Ren Ng

# Light field photography using a handheld plenoptic camera

**Refocusing**

**Novel View Synthesis**

**http://lightfield-forum.com/en/**

# Deep Learning for View Synthesis

# Generating Chairs with CNNs



**Dosovitskiy et al. Learning to Generate Chairs, Tables and Cars with Convolutional Networks**
**PAMI 2017 (CVPR 2015)**

# View Synthesis with Dense Correspondence



**View Synthesis by Appearance Flow**
**Tinghui Zhou, Shubham Tulsiani, Weilun Sun, Jitendra Malik, Alexei A. Efros**
**ECCV 2016**

# View Synthesis with Dense Correspondence



**View Synthesis by Appearance Flow**
**Tinghui Zhou, Shubham Tulsiani, Weilun Sun, Jitendra Malik, Alexei A. Efros**
**ECCV 2016**

# Lots of recent progress using deep learning for view synthesis!

Wiles CVPR 2020

Choi ICCV 2019

Flynn CVPR 2019

The following slides deck is from
Ben Mildenhall*, Pratul Srinivasan*, Matthew Tancik*, Jonathan T. Barron, Ravi Ramamoorthi, Ren Ng

# The problem of novel view interpolation



Inputs: sparsely sampled images of scene

Outputs: *new* views of same scene

# Very successful approach: predict 3D voxel RGB-alpha grid

Input images

Predicted voxel grid

Rendered new views



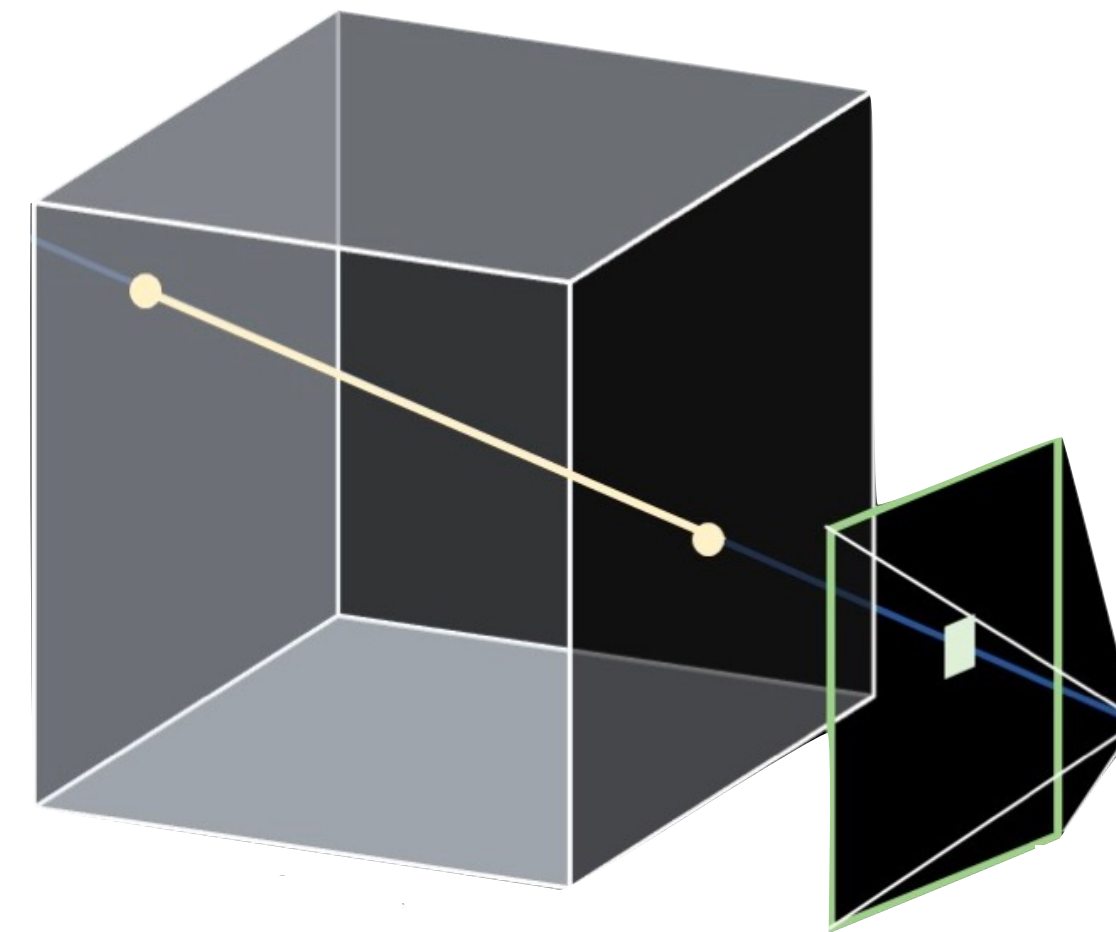Neural Volumes, Lombardi et al. 2019

# RGB-alpha volume rendering for view synthesis

## Soft 3D
(Penner & Zhang 2017)
Culmination of non-deep stereo matching techniques



## 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



Input Sampled View          Promote to MPI

## Neural Volumes
(Lombardi et al. 2019)
Direct gradient descent to optimize an RGBA volume, regularized by a 3D CNN

# RGB-alpha volume rendering for view synthesis

**Soft 3D**
(Penner & Zhang 2017)
Culmination of non-deep stereo matching techniques



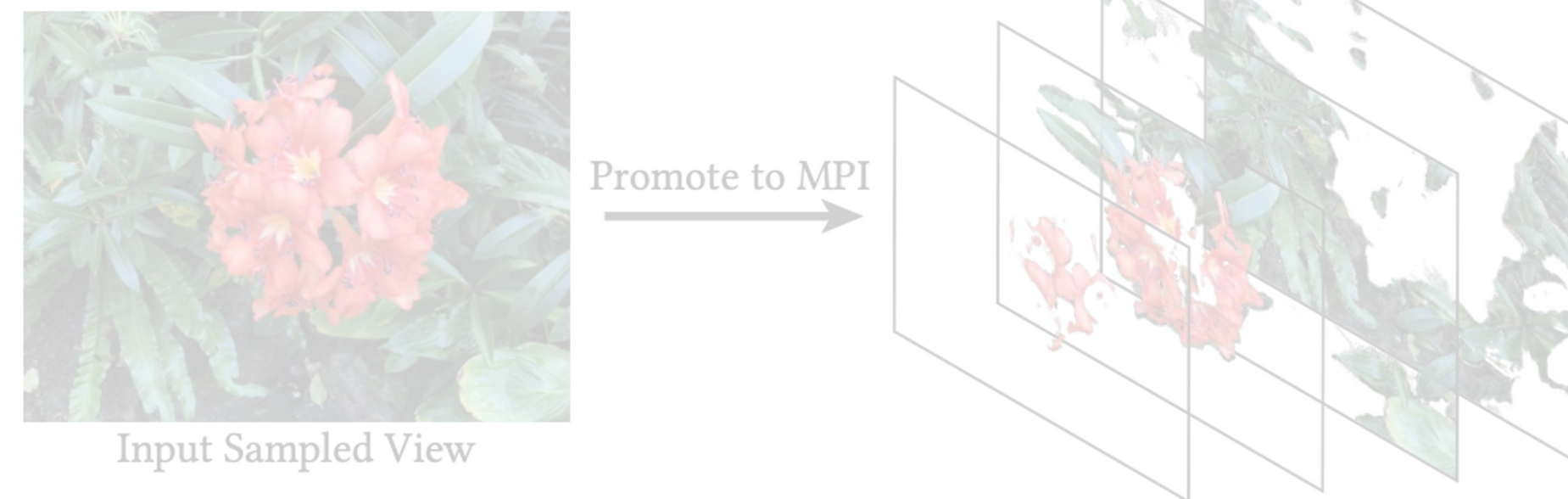**Multiplane image methods**

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



Input Sampled View

Promote to MPI

**Neural Volumes**
(Lombardi et al. 2019)
Direct gradient descent to optimize an RGBA volume, regularized by a 3D CNN
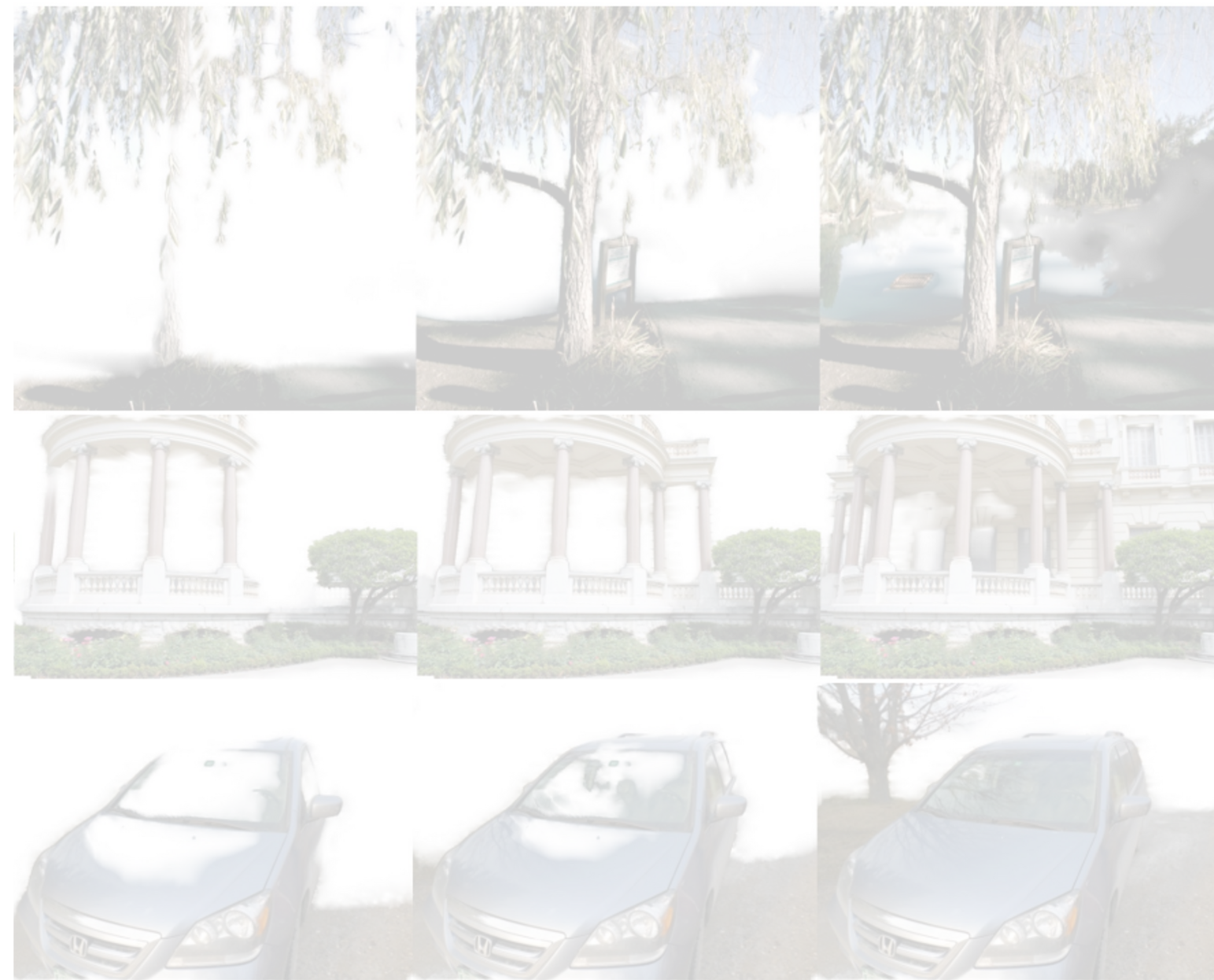
# RGB-alpha volume rendering for view synthesis

**Soft 3D**

(Penner & Zhang 2017)
Culmination of non-deep stereo matching techniques



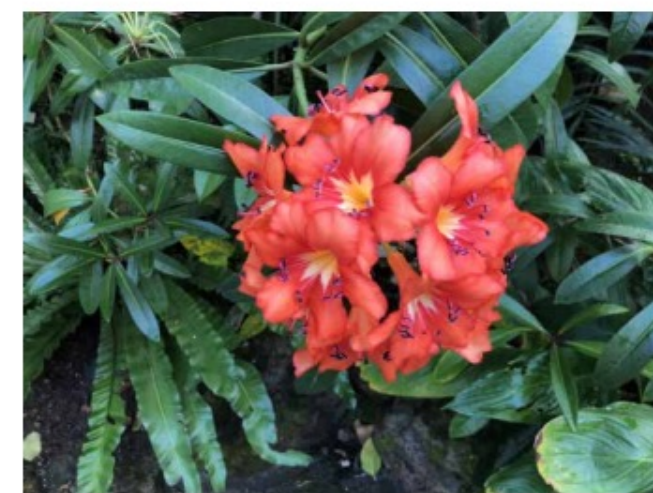**Multiplane image methods**

Stereo Magnification (Zhou et al. 2018)
Pushing the Boundaries… (Srinivasan et al. 2019)
Local Light Field Fusion (Mildenhall et al. 2019)
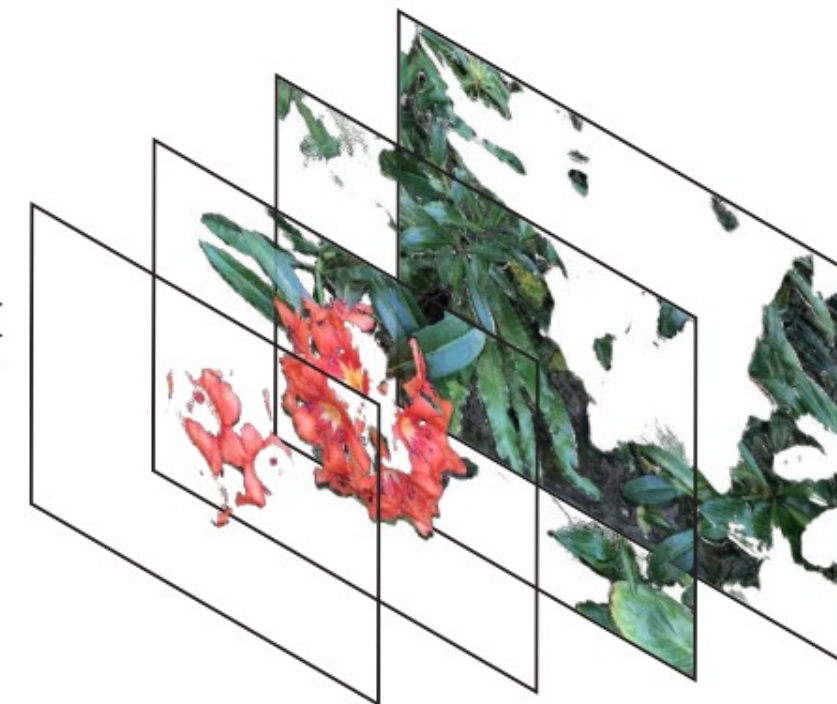DeepView (Flynn et al. 2019)
Single-View… (Tucker & Snavely 2020)

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



Input Sampled View        Promote to MPI

**Neural Volumes**

(Lombardi et al. 2019)
Direct gradient descent to optimize an RGBA volume, regularized by a 3D CNN
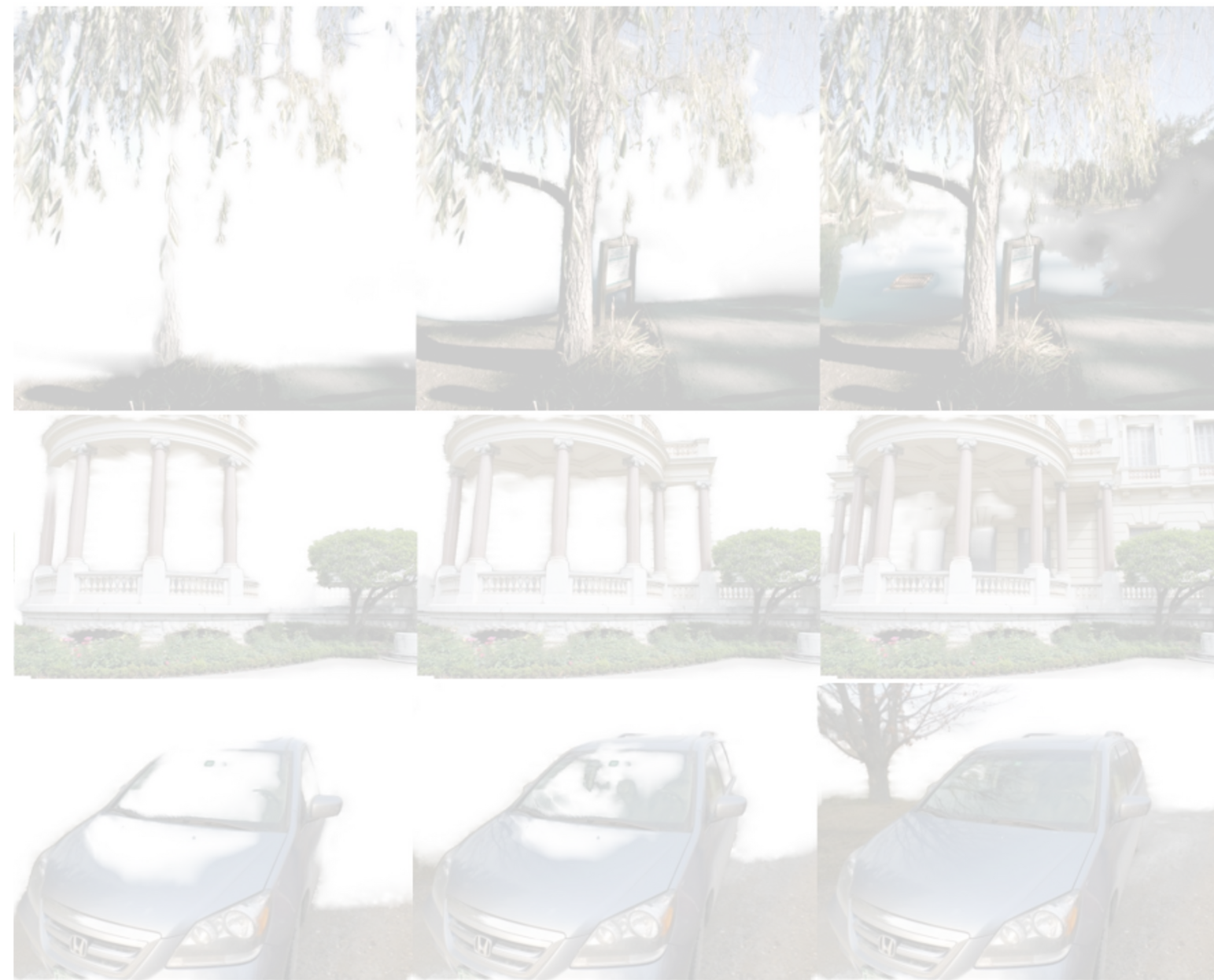


Example Reconstructions
Or a collection of videos from multiple viewpoints
[Lombardi et al. 2019]
Source

# RGB-alpha volume rendering for view synthesis



Input images

Encoder $z$ Decoder

Predicted voxel grid

Rendered new views

Neural Volumes, Lombardi et al. 2019

# RGB-alpha volume rendering for view synthesis

**Soft 3D**
(Penner & Zhang 2017)
Culmination of non-deep stereo matching techniques

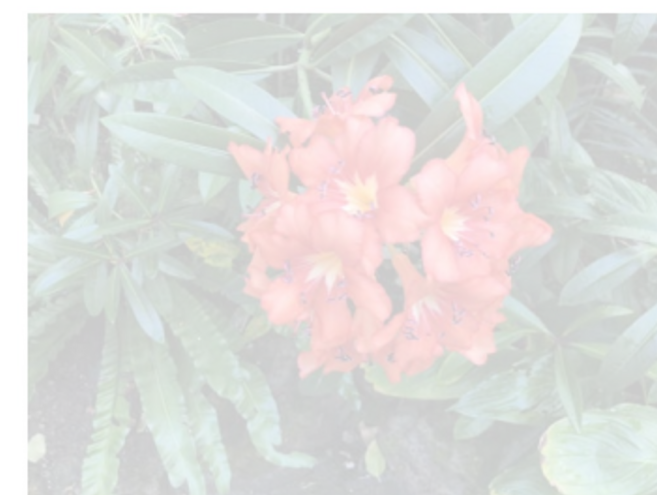**Multiplane image methods**

Stereo Magnification (Zhou et al. 2018)
Pushing the Boundaries… (Srinivasan et al. 2019)
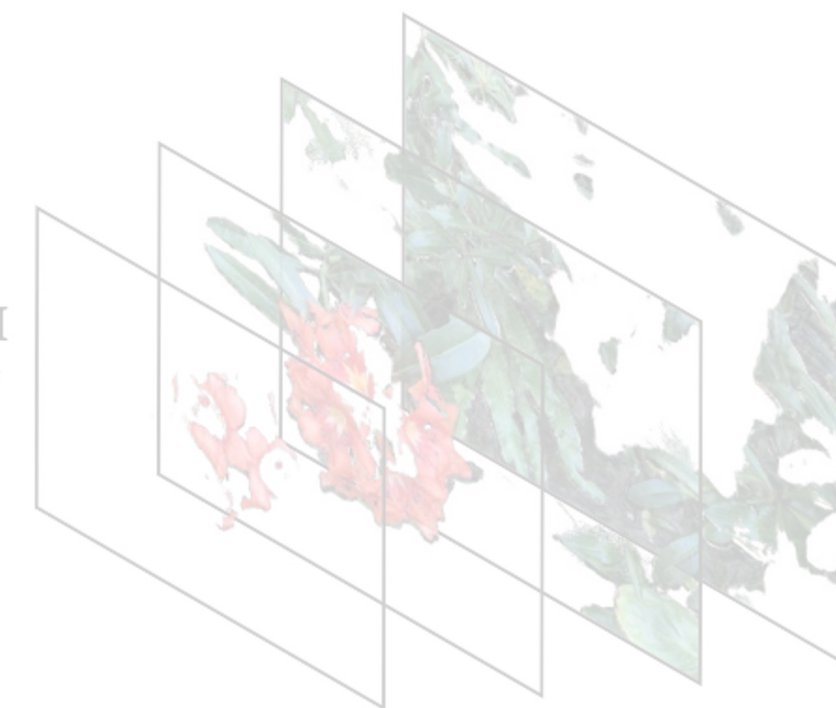Local Light Field Fusion (Mildenhall et al. 2019)
DeepView (Flynn et al. 2019)
Single-View… (Tucker & Snavely 2020)

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

**Neural Volumes**
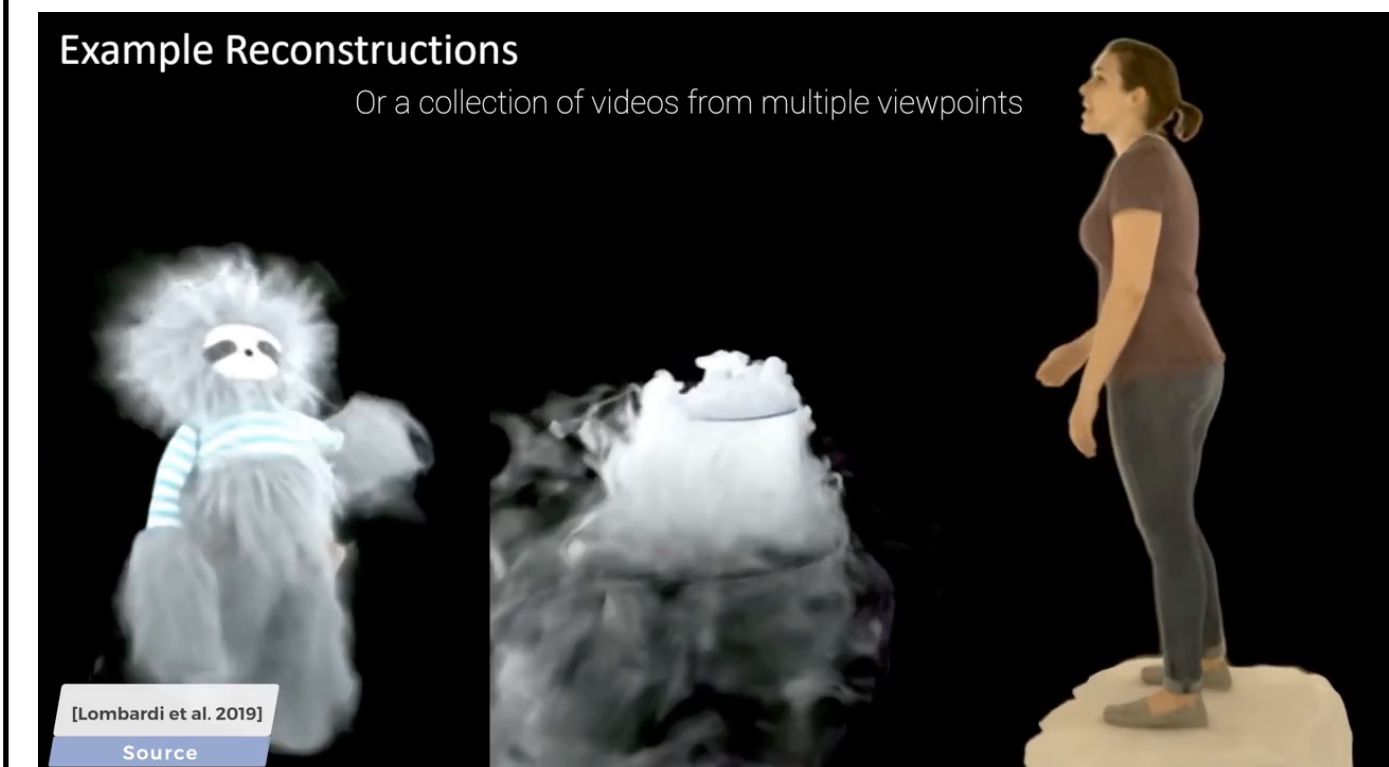(Lombardi et al. 2019)
Direct gradient descent to optimize an RGBA volume, regularized by a 3D CNN

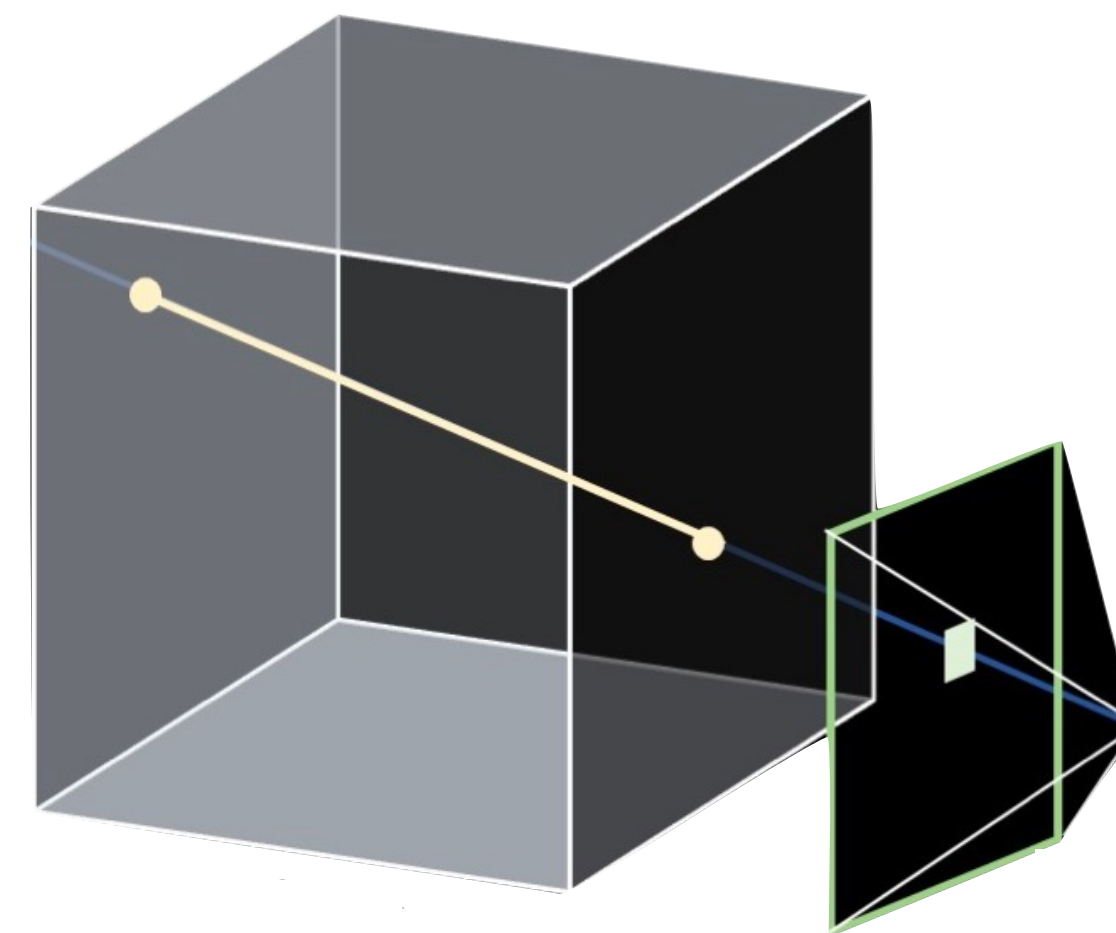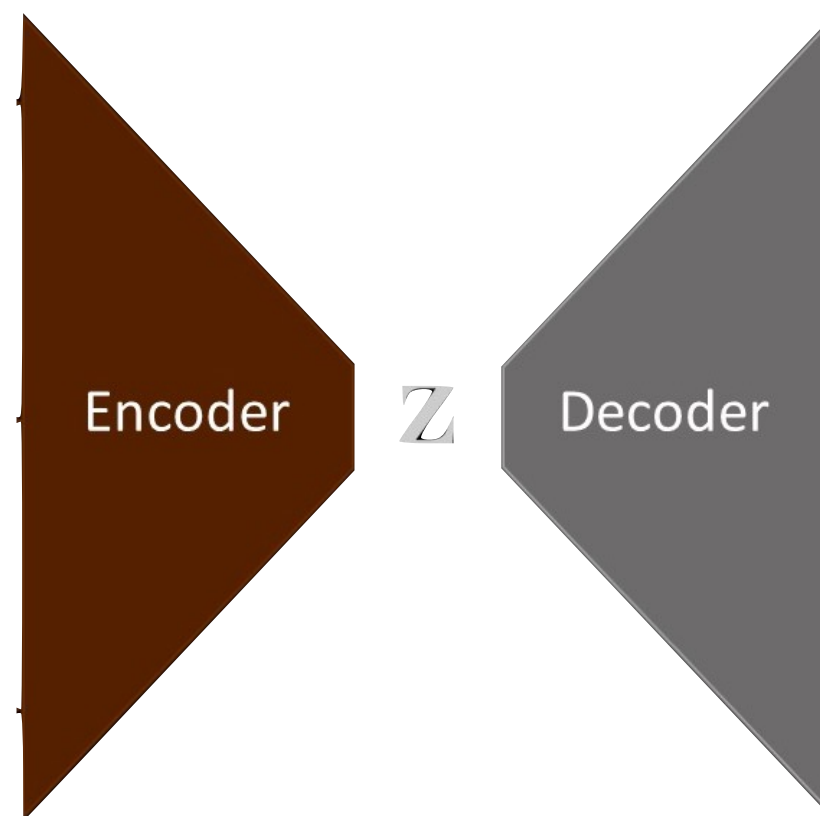+ Great rendering model: good for optimization
- Horrible storage requirements (1-10 GB)

# Neural networks as a continuous shape representation



$$f_\theta(p) = \tau$$

Occupancy Networks, Mescheder et al. CVPR 2019

43

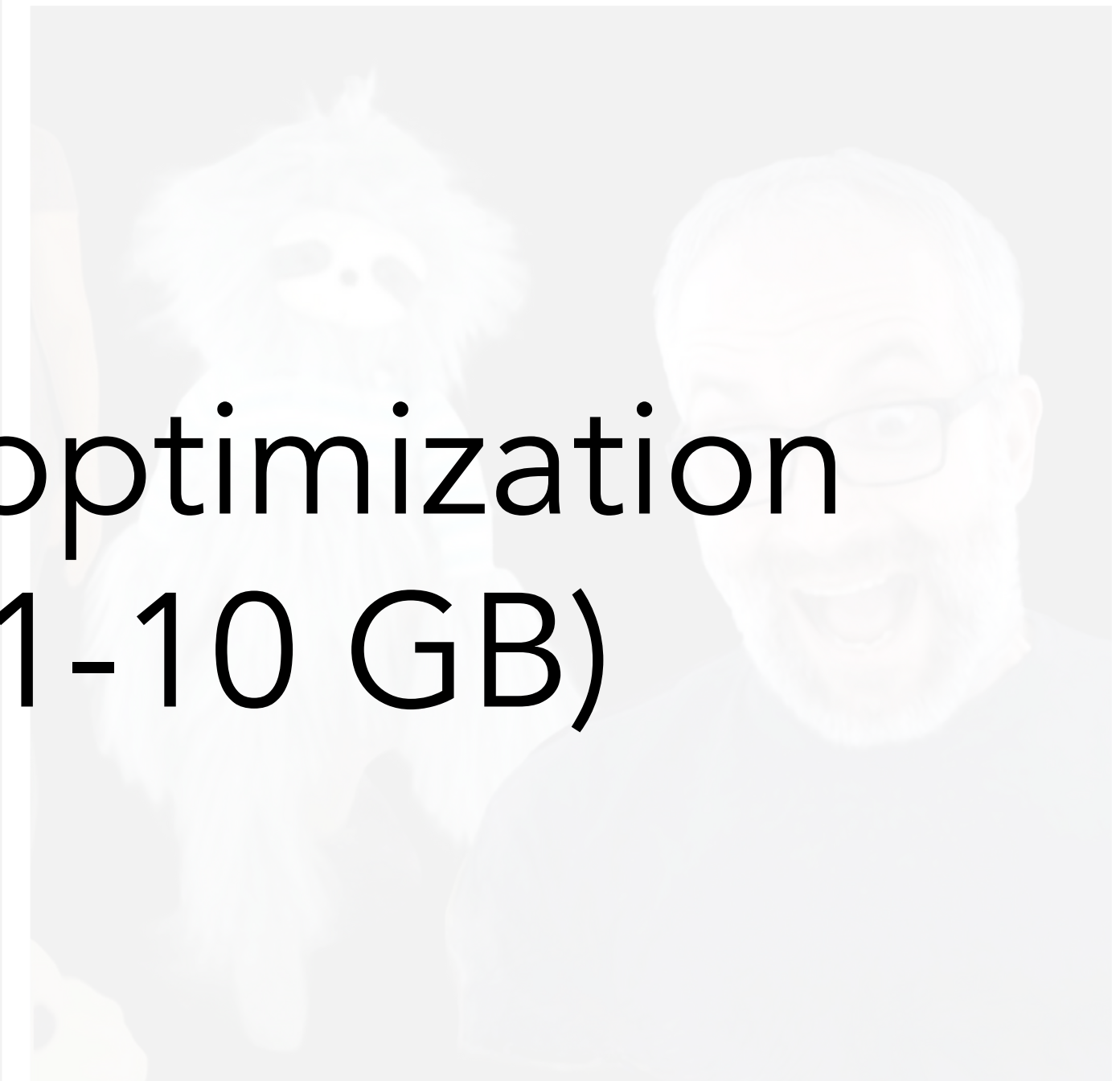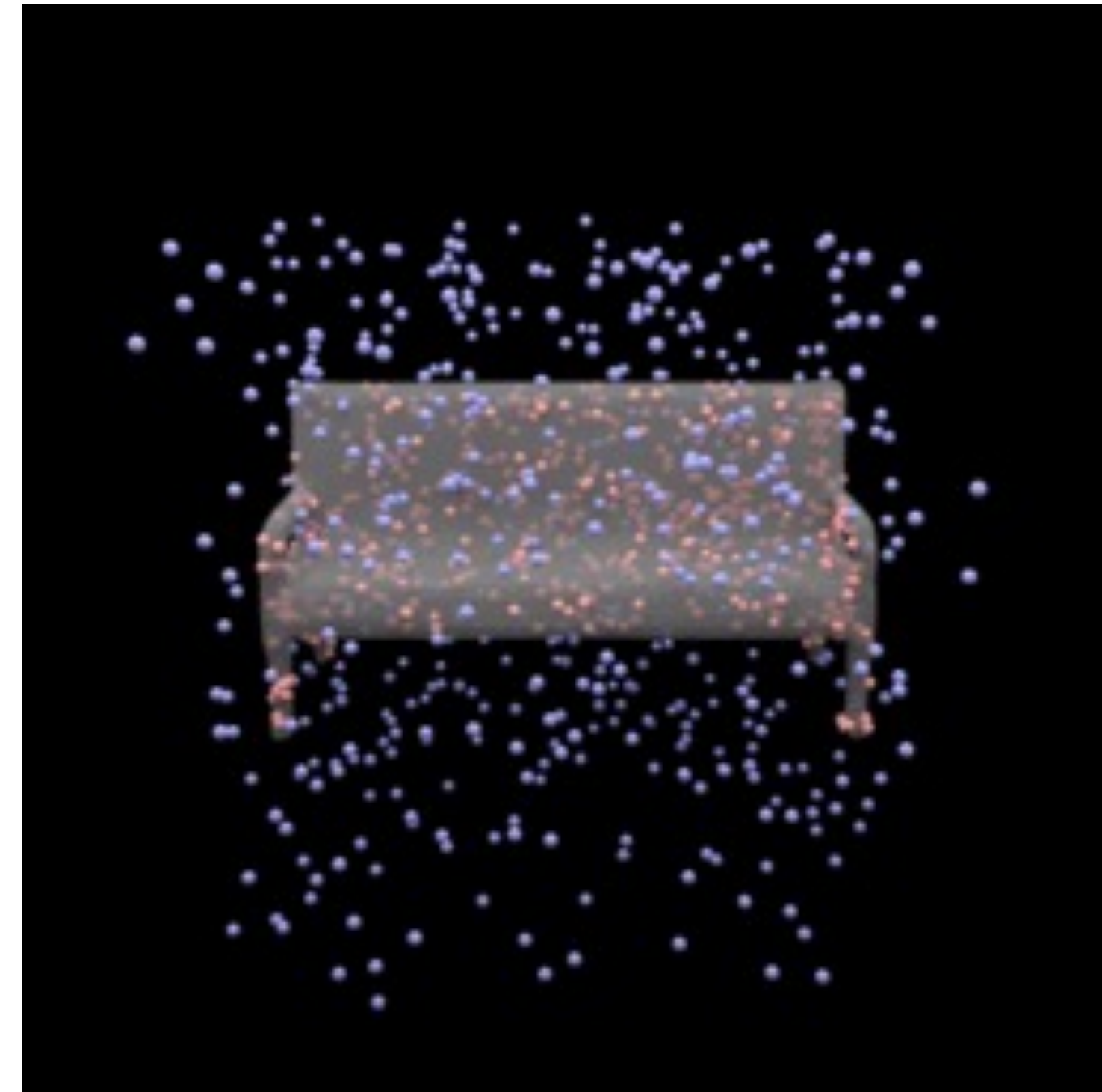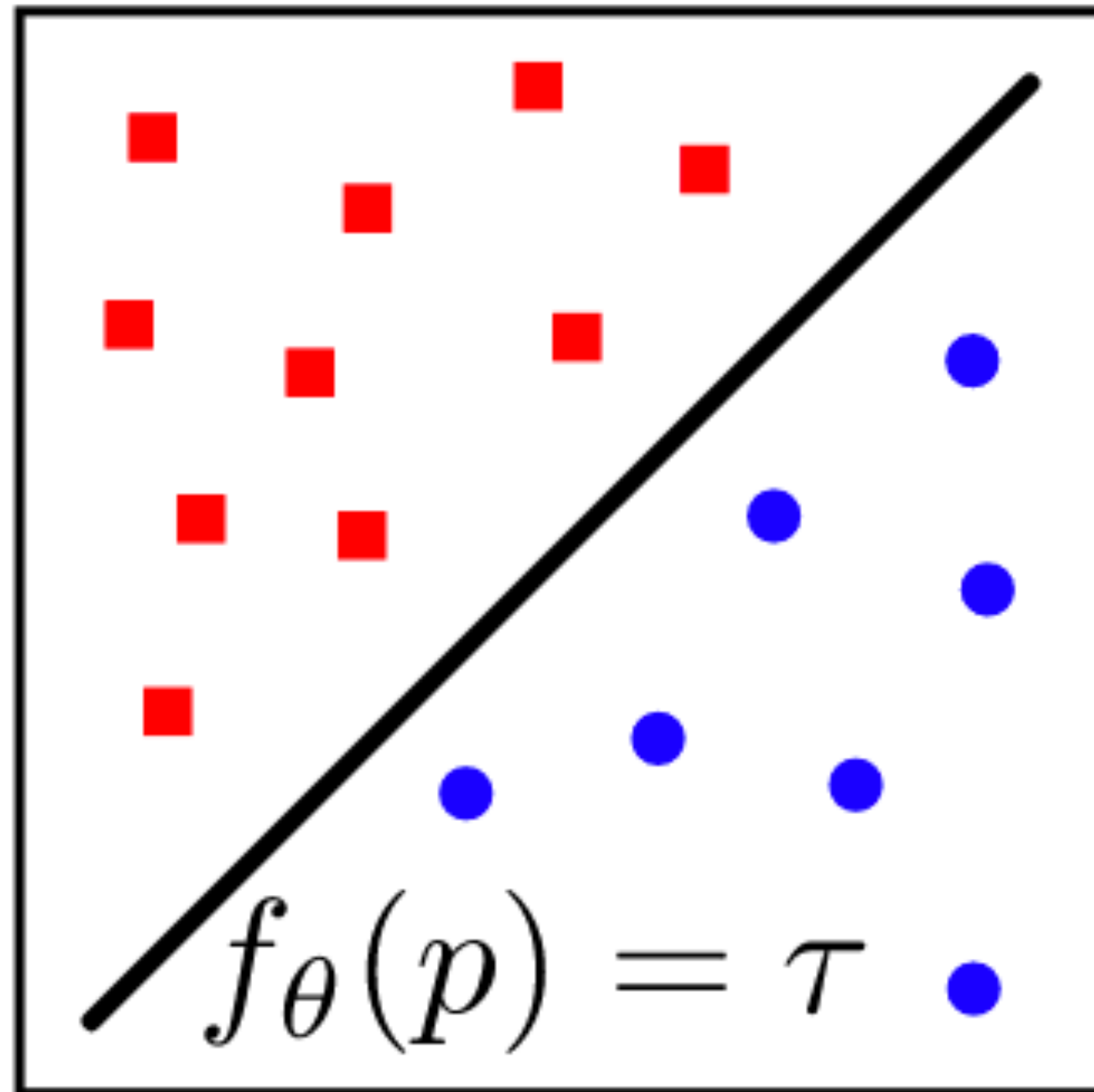# Neural networks as a continuous shape representation

**Occupancy Networks**
(Mescheder et al. 2019)
$(x, y, z) \rightarrow occupancy$



$$f_\theta(p) = \tau$$

**DeepSDF**
(Park et al. 2019)
$(x, y, z) \rightarrow distance$



**Scene Representation Networks**
(Sitzmann et al. 2019)
$(x, y, z) \rightarrow latent\ vec.\ (color,\ dist.)$



**Differentiable Volumetric Rendering**
(Niemeyer et al. 2020)
$(x, y, z) \rightarrow color,\ occ.$

$$\frac{\partial \hat{d}}{\partial \theta} = -\left( \frac{\partial f_\theta(\hat{\mathbf{p}})}{\partial \hat{\mathbf{p}}} \cdot \mathbf{w} \right)^{-1} \frac{\partial f_\theta(\hat{\mathbf{p}})}{\partial \theta}$$

# Neural networks as a shape representation

**DeepSDF**
(Park et al. 2019)
$(x, y, z) \rightarrow distance$

**Occupancy Networks**
(Mescheder et al. 2019)
$(x, y, z) \rightarrow occupancy$

## - Limited rendering model: difficult to optimize
## + Highly compressible (1-10 MB)

**Scene Representation Networks**
(Sitzmann et al. 2019)
$(x, y, z) \rightarrow latent\ vec.\ (color,\ dist.)$

**Differentiable Volumetric Rendering**
(Niemeyer et al. 2020)
$(x, y, z) \rightarrow color,\ occ.$

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

# NeRF achieves state-of-the-art results on an extremely difficult problem



LLFF [Mildenhall 2019]

NeRF

Nearest Input

# NeRF achieves state-of-the-art results on an extremely difficult problem



Neural Volumes [Lombardi 2019]

NeRF

Nearest Input

# NeRF achieves state-of-the-art results on an extremely difficult problem



SRN [Sitzmann 2019]

NeRF

Nearest Input

# Key points

‣ Continuous neural network as a volumetric scene representation (5D = xyz + direction)

‣ Use volume rendering model to synthesize new views

‣ Optimize using rendering loss for one scene (no prior training)

‣ One extra trick for passing coordinates into network to get high frequency details

# Key points

‣ **Continuous neural network as a volumetric scene representation (5D = xyz + direction)**

‣ Use volume rendering model to synthesize new views

‣ Optimize using rendering loss for one scene (no prior training)

‣ One extra trick for passing coordinates into network to get high frequency details

# Representing a scene as a continuous 5D function

$$(x, y, z, \theta, \phi)$$

$$\underbrace{x, y, z}_{\text{Spatial location}}, \underbrace{\theta, \phi}_{\text{Viewing direction}}$$

$$F_\Omega$$

Fully-connected
neural network
9 layers,
256 channels

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

$$\underbrace{r, g, b}_{\text{Output color}}, \underbrace{\sigma}_{\text{Output density}}$$

# Neural network replaces large N-d array



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

$$(x, y, z, \theta, \phi) \longrightarrow \boxed{\phantom{F}} \longrightarrow (r, g, b, \sigma)$$
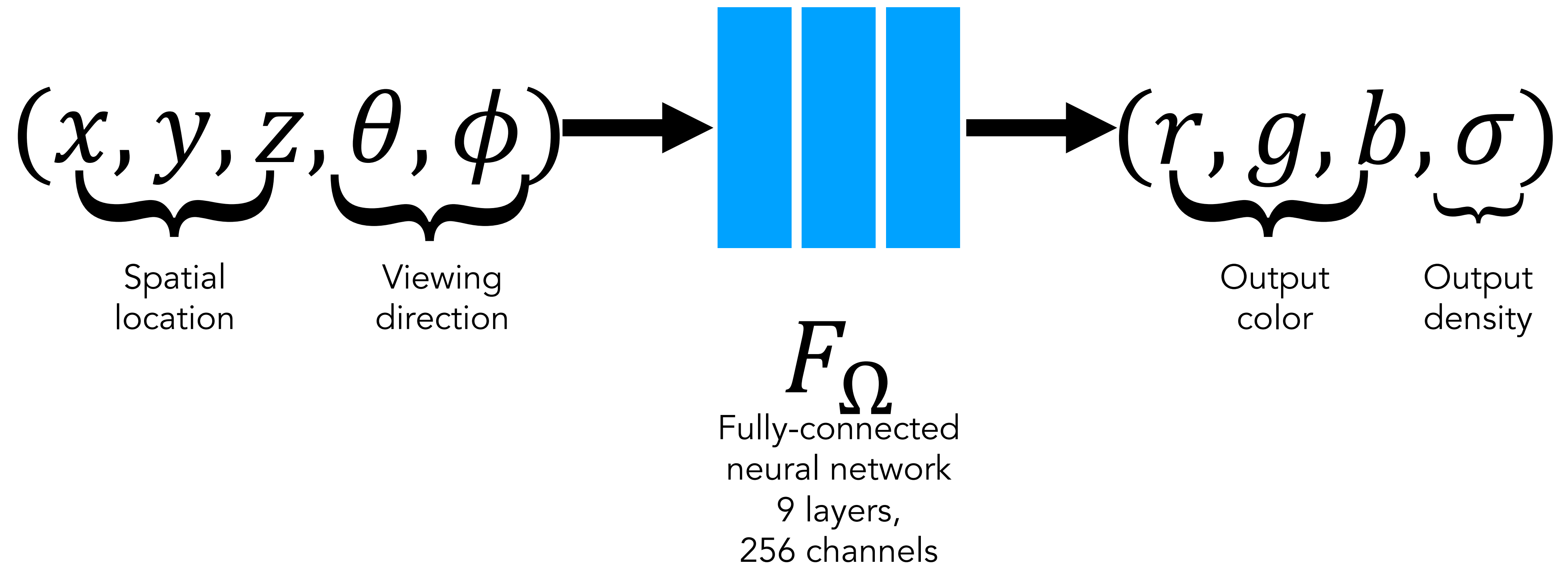
$F_\Omega$

# Key points

‣ Continuous neural network as a volumetric scene representation (5D = xyz + direction)

‣ **Use volume rendering model to synthesize new views**

‣ Optimize using rendering loss for one scene (no prior training)

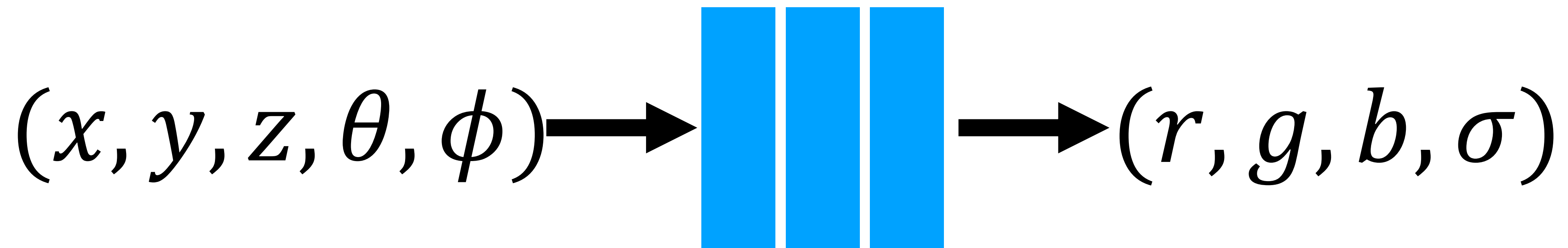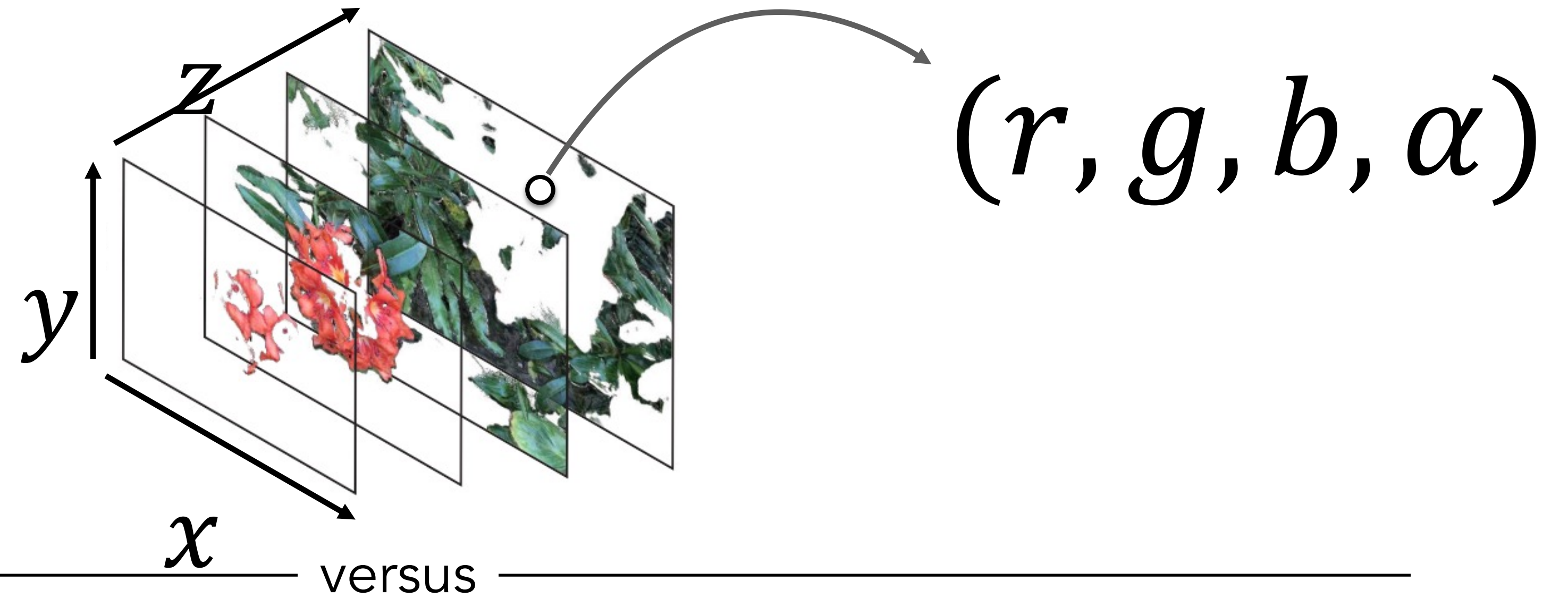‣ One extra trick for passing coordinates into network to get high frequency details

# Generate views with traditional volume rendering

# Generate views with traditional volume rendering

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



Ray

$t_N$

3D volume

$t_1$

Camera

# 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



Ray

$t_N$

3D volume

$t_1$

Camera

# Generate views with traditional volume rendering

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

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

colors

weights

How much light is blocked earlier along ray:

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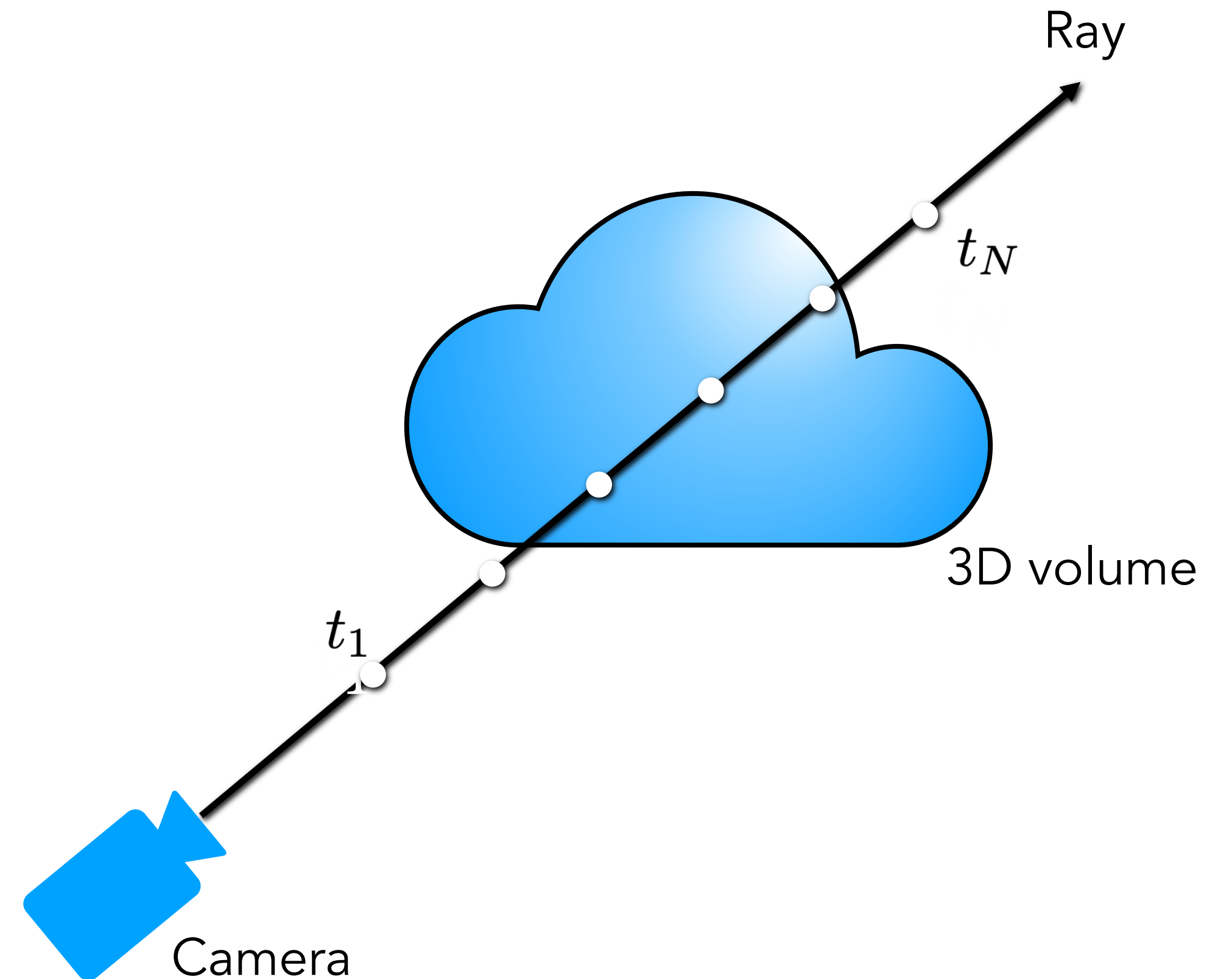
Ray

$t_N$

colors

$t_1$

$T_i$

3D volume

Camera

# Generate views with traditional volume rendering

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

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

colors

weights

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 $i$:

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

Ray

$t_N$

$\alpha_i$

3D volume

$t_1$

$T_i$

Camera

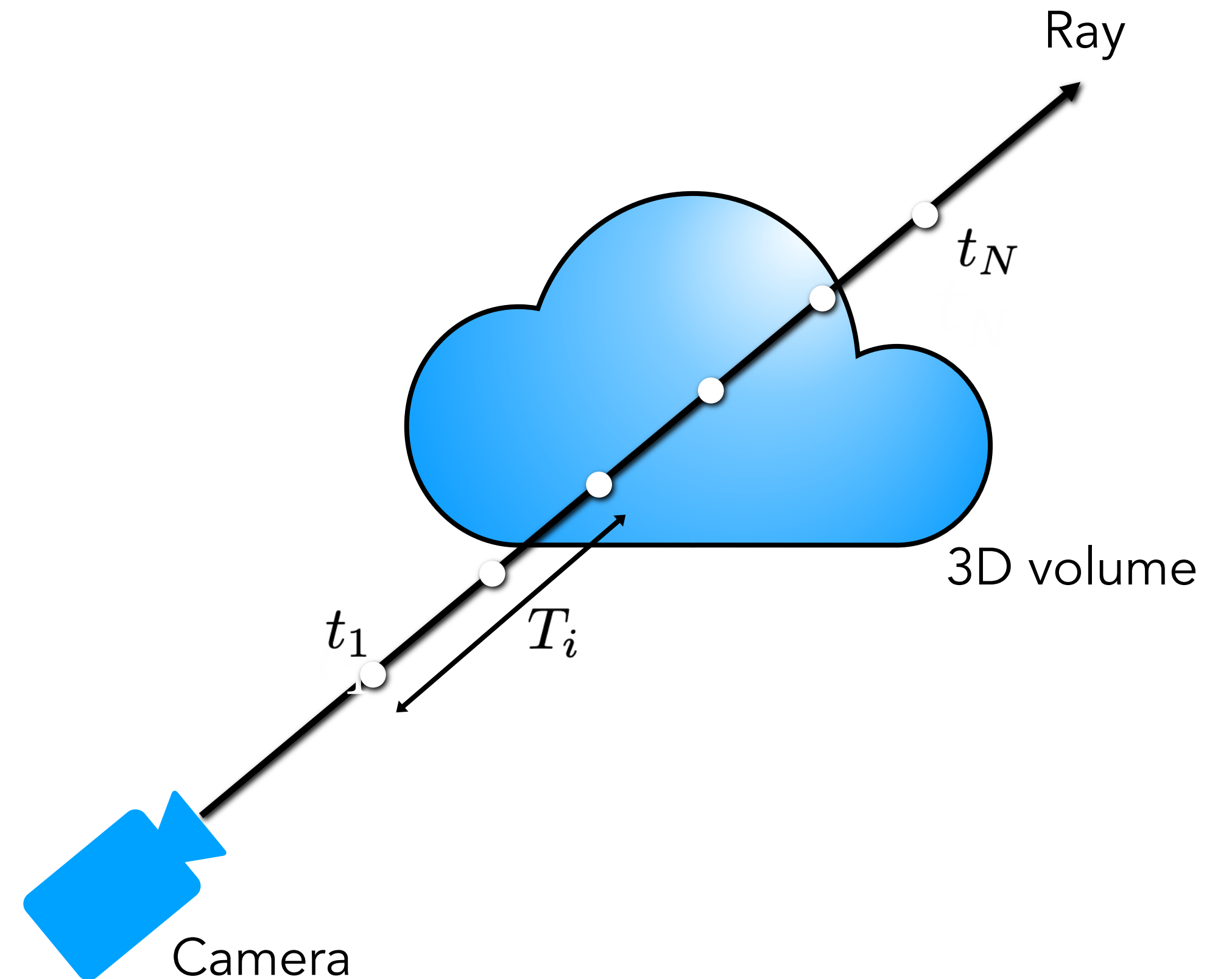# Sigma parametrization for continuous opacity

Rendering model for ray r(t) = o + td:
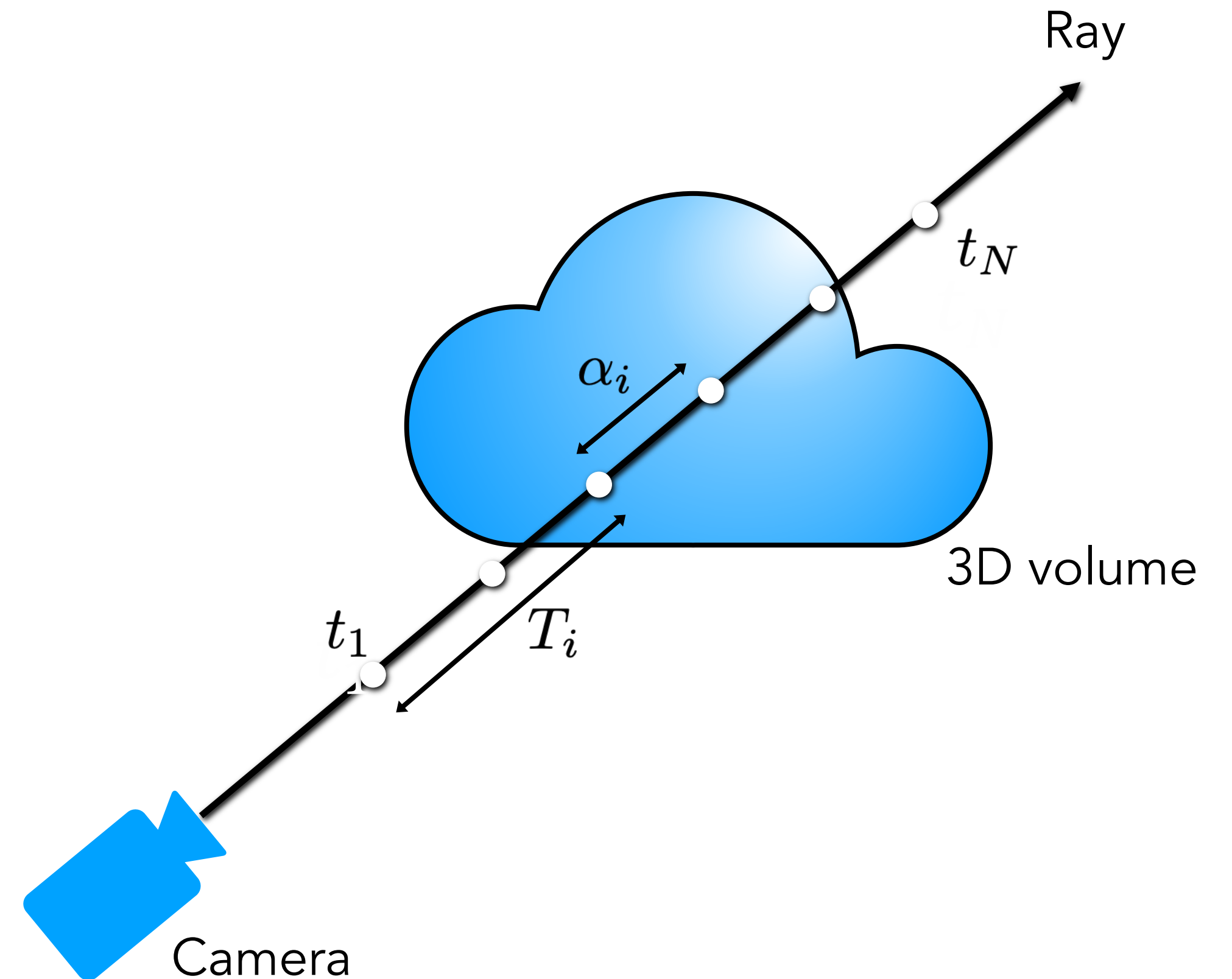
$$C \approx \sum_{i=1}^{N} T_i \alpha_i c_i$$

colors

weights

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 $i$:

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

Ray

$t_N$

3D volume

$t_1$

Camera

# Effective resolution is tied to distance between samples

Rendering model for ray r(t) = o + td:
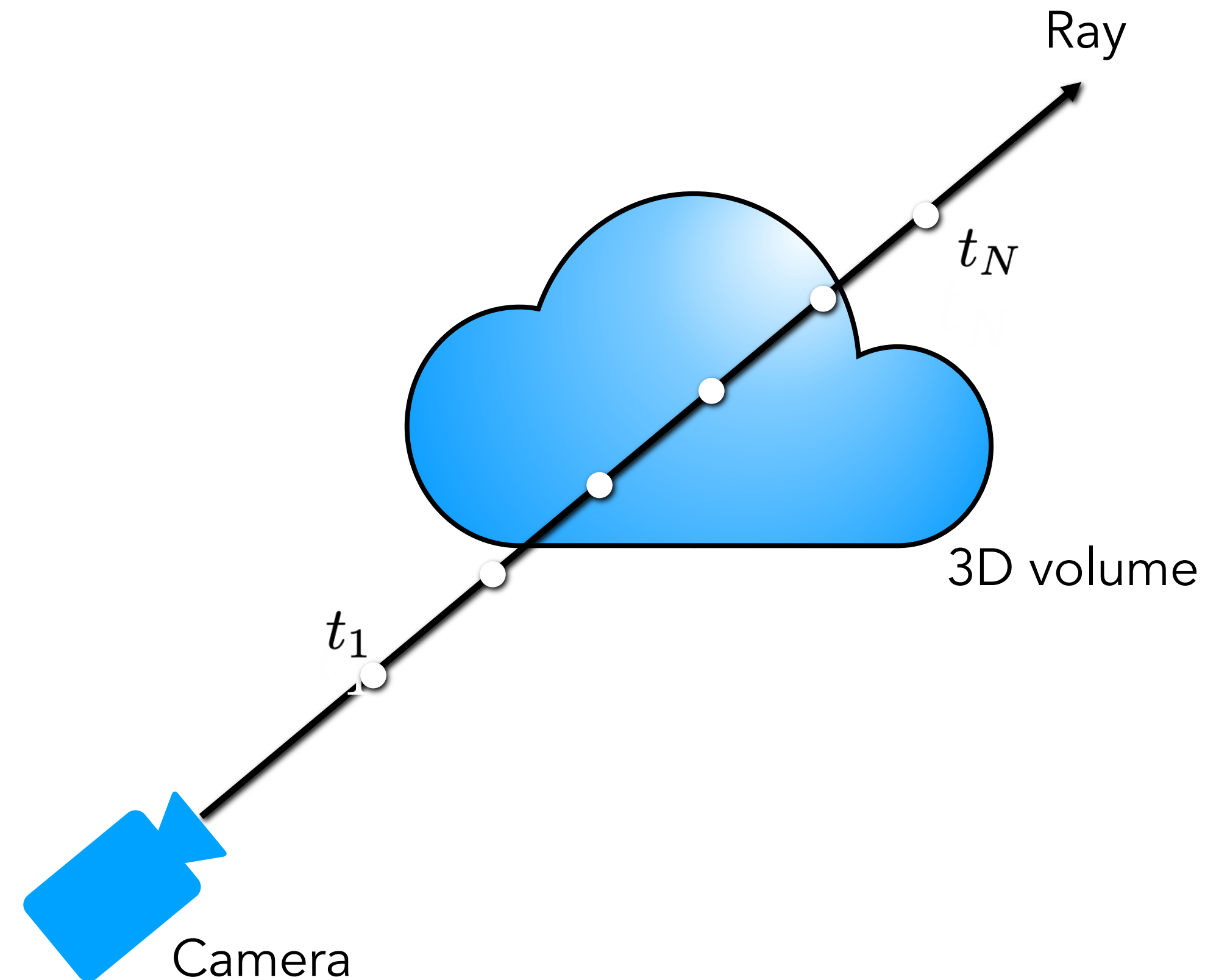
$$C \approx \sum_{i=1}^{N} T_i \alpha_i c_i$$
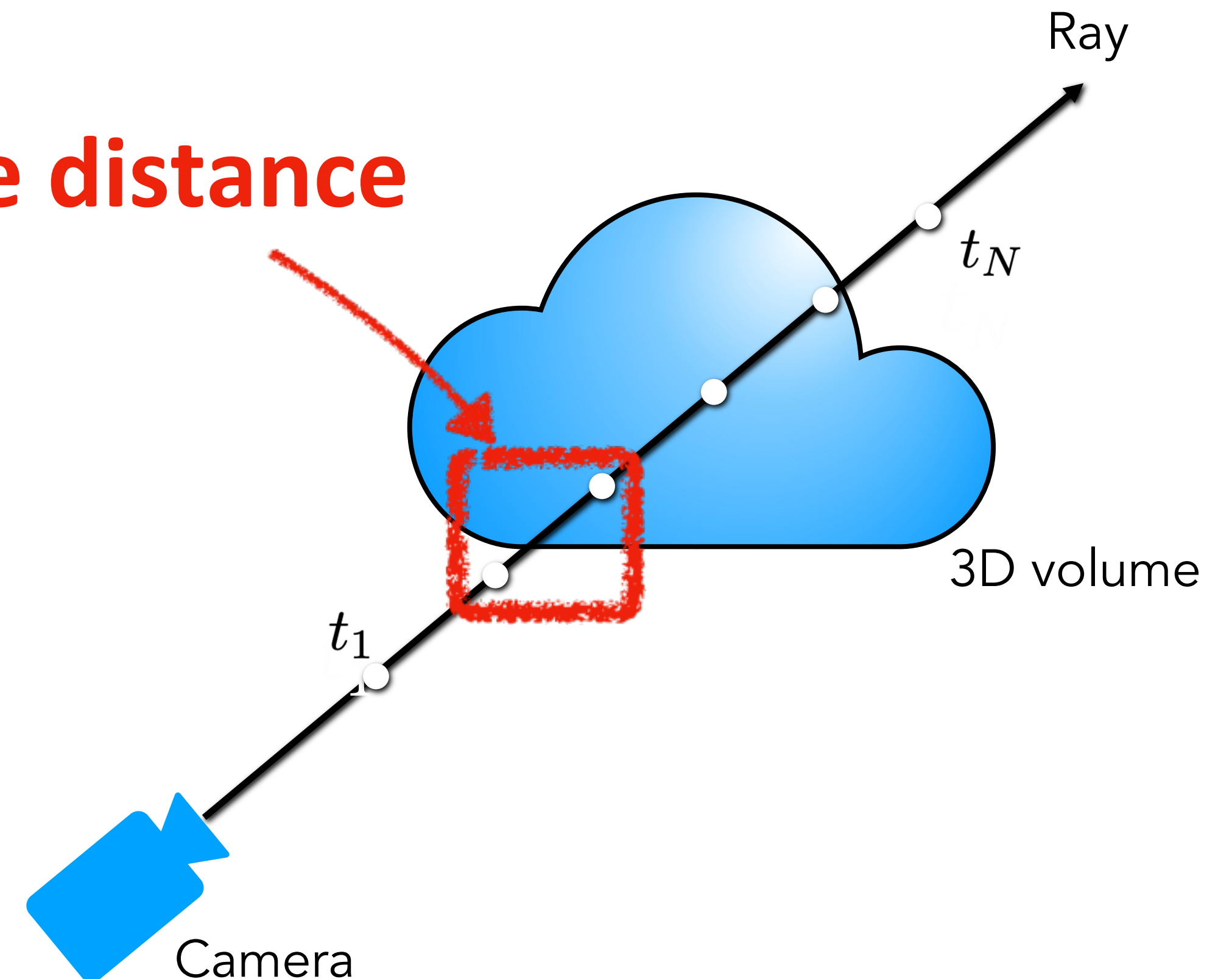
colors

weights

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 $i$:

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

Ray

**sample distance**

$t_N$

3D volume

$t_1$

Camera

# Can we allocate samples more efficiently? Two pass rendering



Ray

3D volume

Camera

# Two pass rendering: coarse

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

treat weights as probability
distribution for new samples

Ray

3D volume

Camera

# Two pass rendering: fine

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

treat weights as probability
distribution for new samples

Ray

3D volume

Camera

# Viewing directions as input



Ray

$(x, y, z, \theta, \phi)$ as input

3D volume

Camera

# Viewing directions as input

Ray

Manipulate $(\theta, \phi)$ to visualize
view-dependent effects

3D volume

Camera

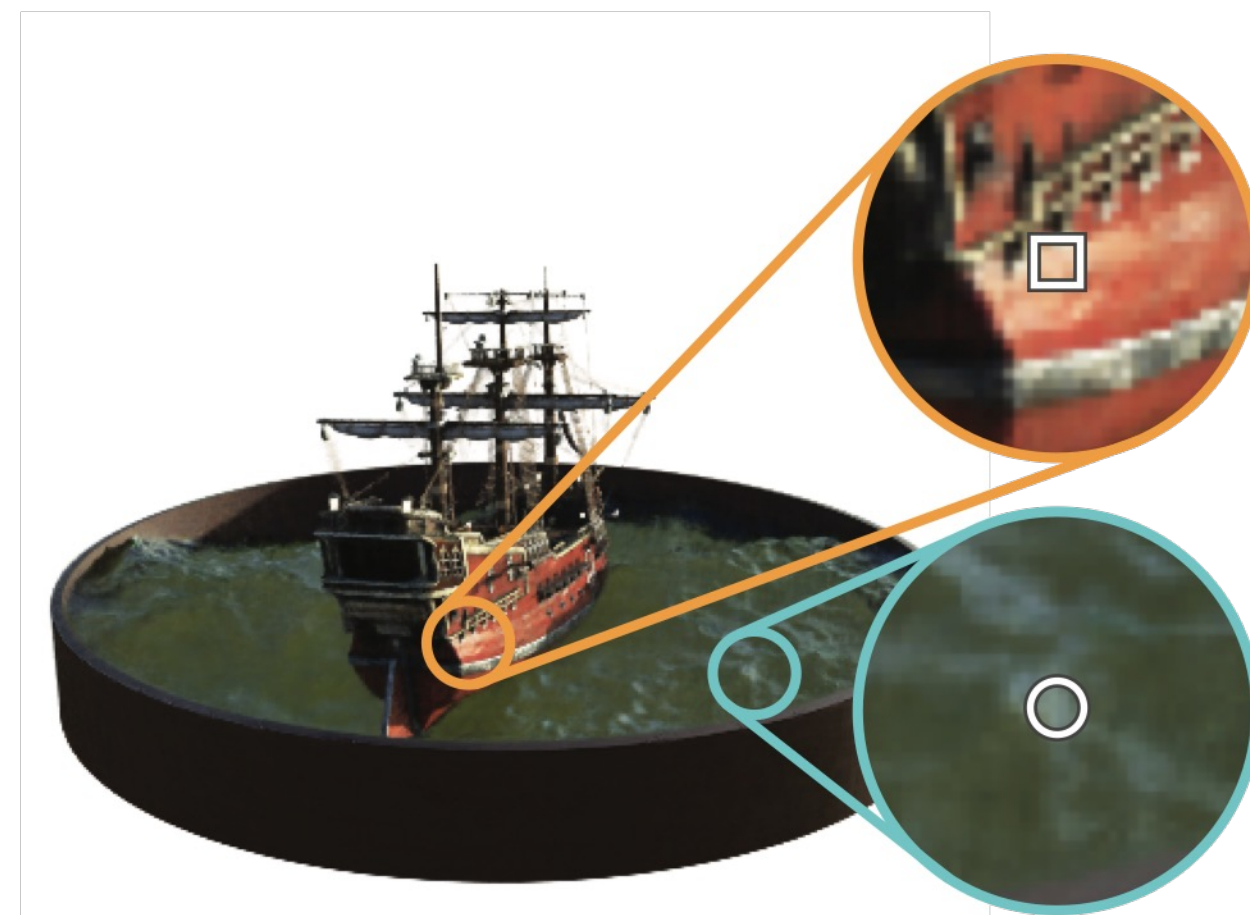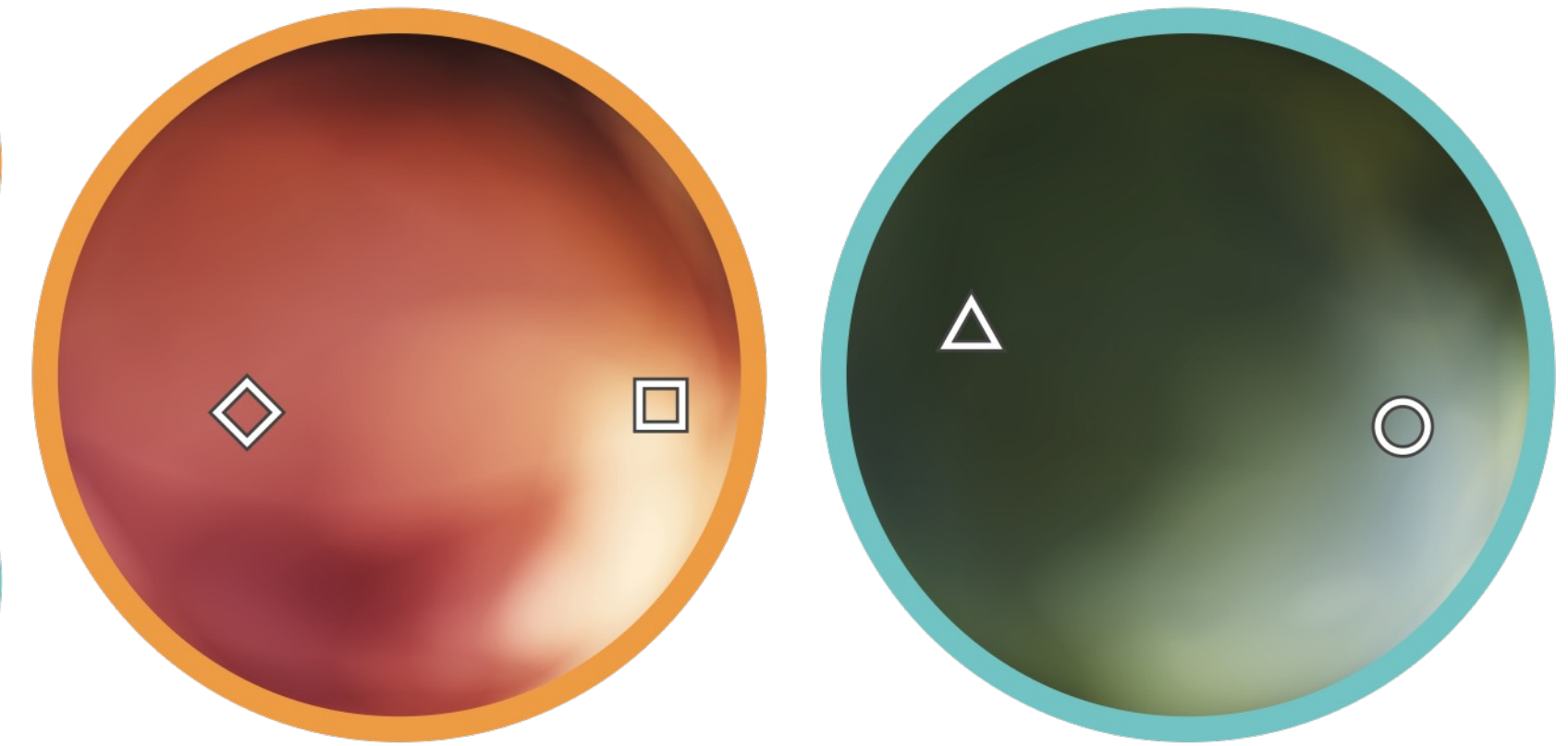# Viewing directions as input



(a) View 1    (b) View 2    (c) Radiance Distributions

# Key points

‣ Continuous neural network as a volumetric scene representation (5D = xyz + direction)

‣ Use volume rendering model to synthesize new views

‣ **Optimize using rendering loss for one scene (no prior training)**

‣ One extra trick for passing coordinates into network to get high frequency details

# Volume rendering is trivially differentiable

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

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

**differentiable w.r.t.**

colors

weights

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 $i$:

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

Ray

$t_N$

$\alpha_i$

3D volume

$t_1$  $T_i$

Camera

# Optimize with gradient descent on rendering loss



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

# Training network to reproduce all input views of the scene

# Naive implementation produces blurry results



NeRF (Naive)

# Naive implementation produces blurry results



NeRF (Naive)                    NeRF (with positional encoding)

# Key points

‣ Continuous neural network as a volumetric scene representation (5D = xyz + direction)

‣ Use volume rendering model to synthesize new views

‣ Optimize using rendering loss for one scene (no prior training)

‣ **One extra trick for passing coordinates into network to get high frequency details**

# Challenge:

How to get MLPs to represent higher frequency functions?

# Simpler toy problem: memorizing a 2D image

$$(x, y) \longrightarrow \blacksquare \longrightarrow (r, g, b)$$
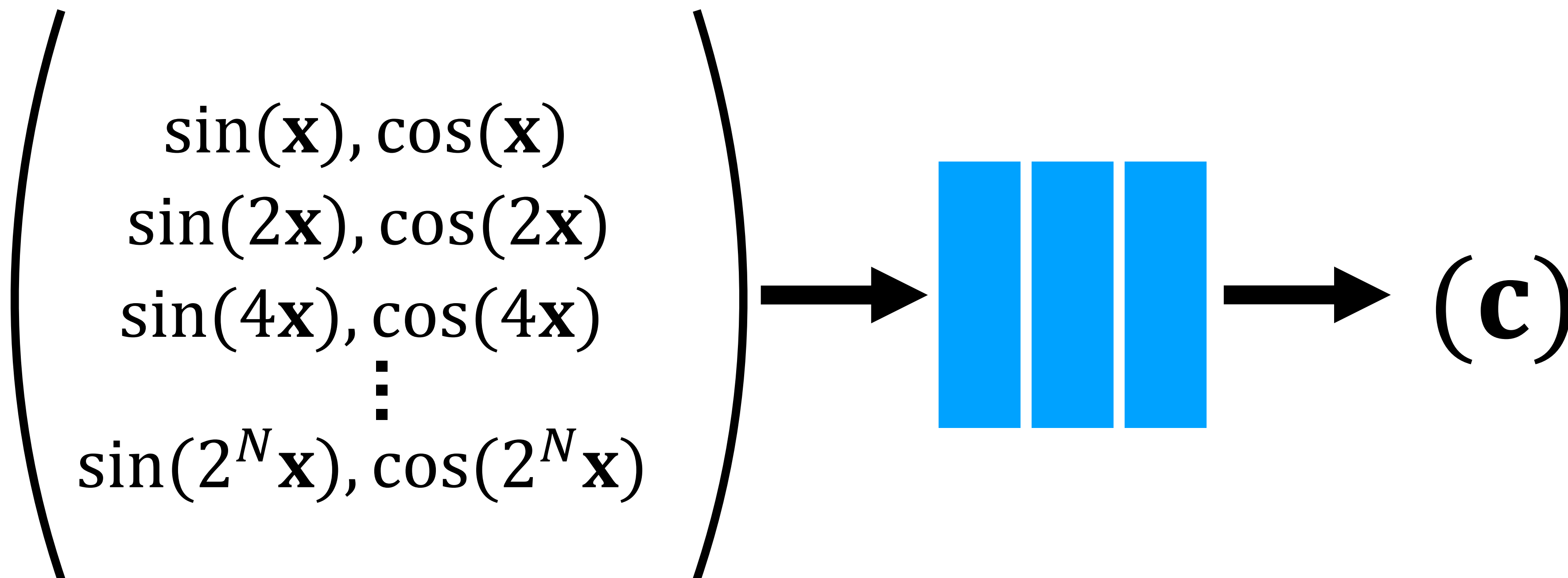
# Simple trick enables network to memorize images

Ground truth image

Standard fully-connected net

# Positional encoding: high frequency embedding of input coordinates

$$(\mathbf{x}) \longrightarrow \blacksquare\blacksquare\blacksquare \longrightarrow (\mathbf{c})$$

$$\begin{pmatrix} \sin(\mathbf{x}), \cos(\mathbf{x}) \\ \sin(2\mathbf{x}), \cos(2\mathbf{x}) \\ \sin(4\mathbf{x}), \cos(4\mathbf{x}) \\ \vdots \\ \sin(2^N\mathbf{x}), \cos(2^N\mathbf{x}) \end{pmatrix} \longrightarrow \blacksquare\blacksquare\blacksquare \longrightarrow (\mathbf{c})$$

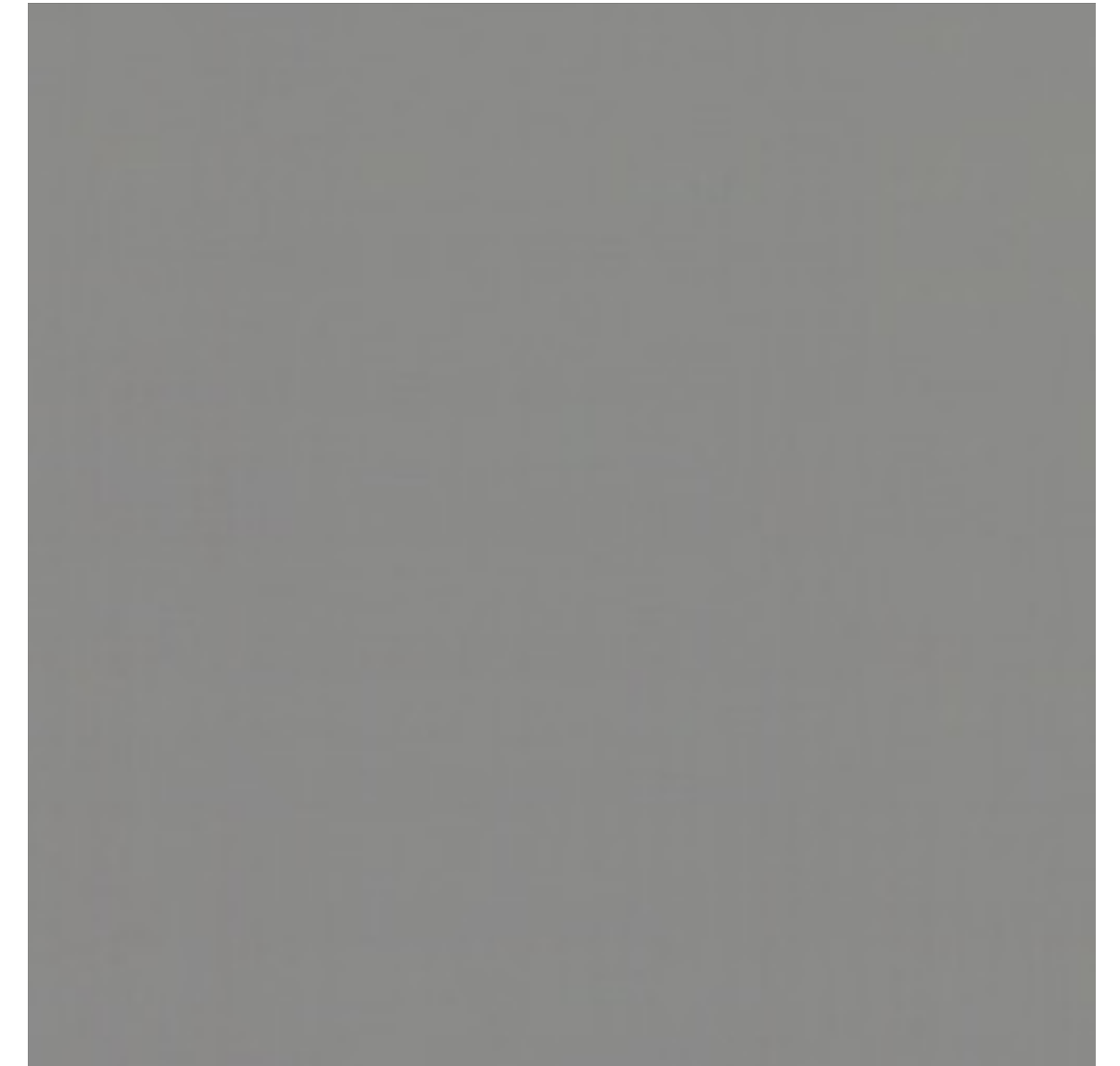# Simple trick enables network to memorize images

Ground truth image

Standard fully-connected net

With "embedding"

# Positional encoding also directly improves our scene representation!



NeRF (Naive)
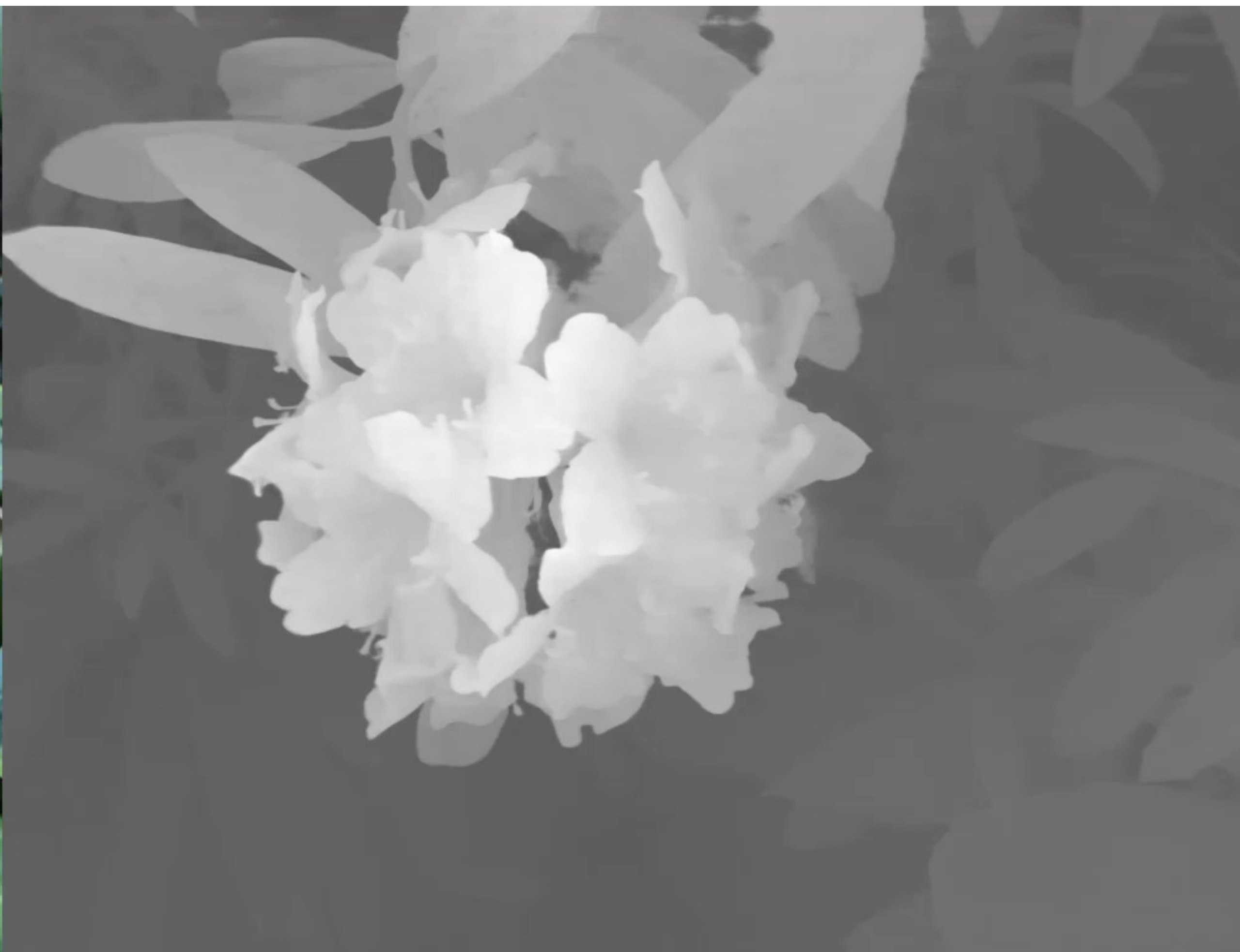
NeRF (with positional encoding)

# Results

# NeRF encodes convincing view-dependent effects using directional dependence

# NeRF encodes convincing view-dependent effects using directional dependence
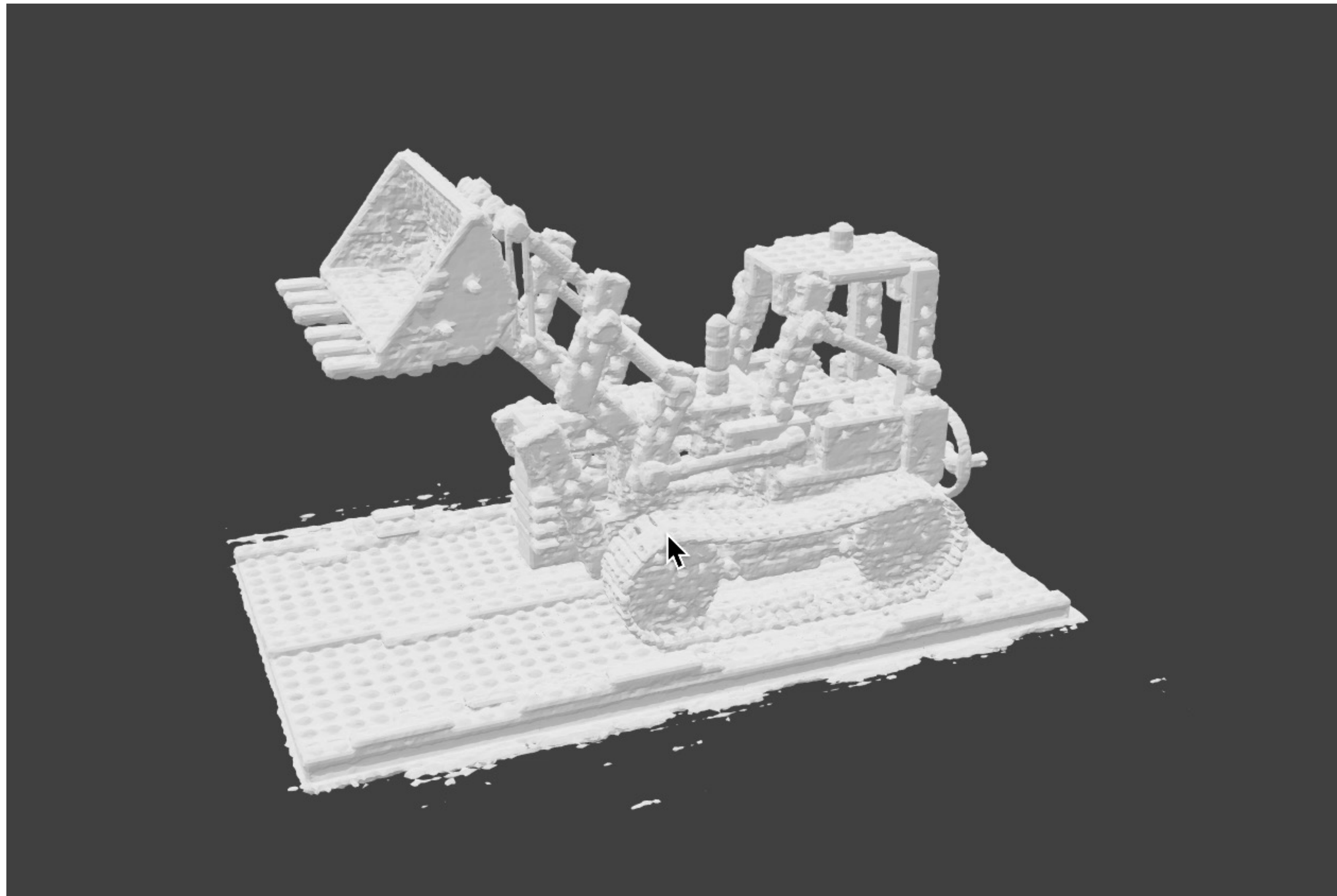
# NeRF encodes detailed scene geometry with occlusion effects

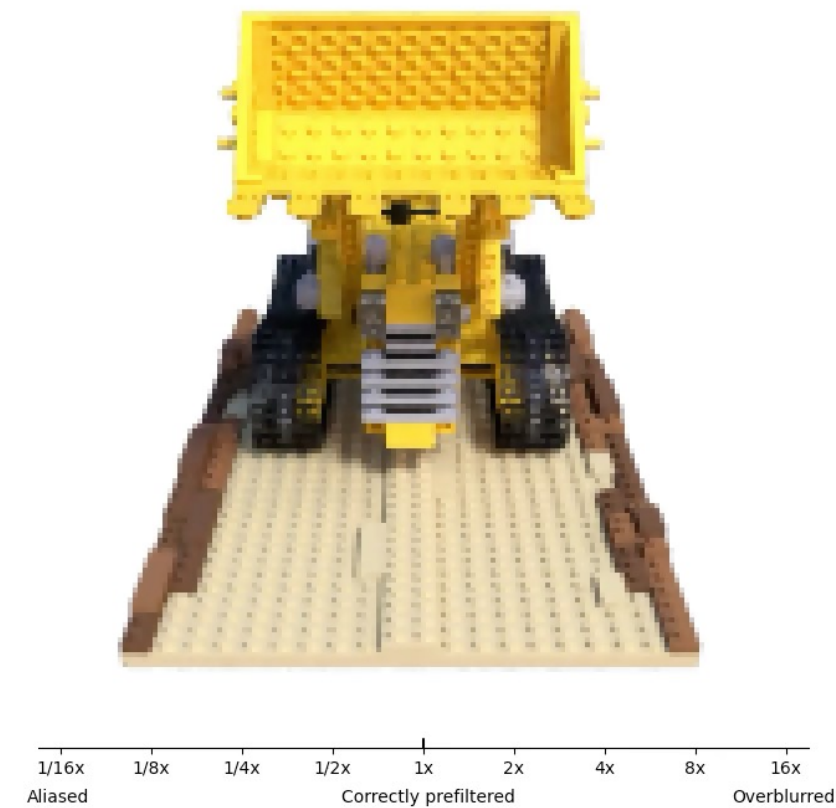# NeRF encodes detailed scene geometry with occlusion effects

# NeRF encodes detailed scene geometry with occlusion effects

# NeRF encodes detailed scene geometry

# Thank You!



16-726, Spring 2023
https://learning-image-synthesis.github.io/sp23