# 3D Perception: PointNet and NERFs

Sanjiban Choudhury

# Last Class: How does a robot identify objects?

# Last Class: How does a robot identify objects?

**Classification**



**CAT**

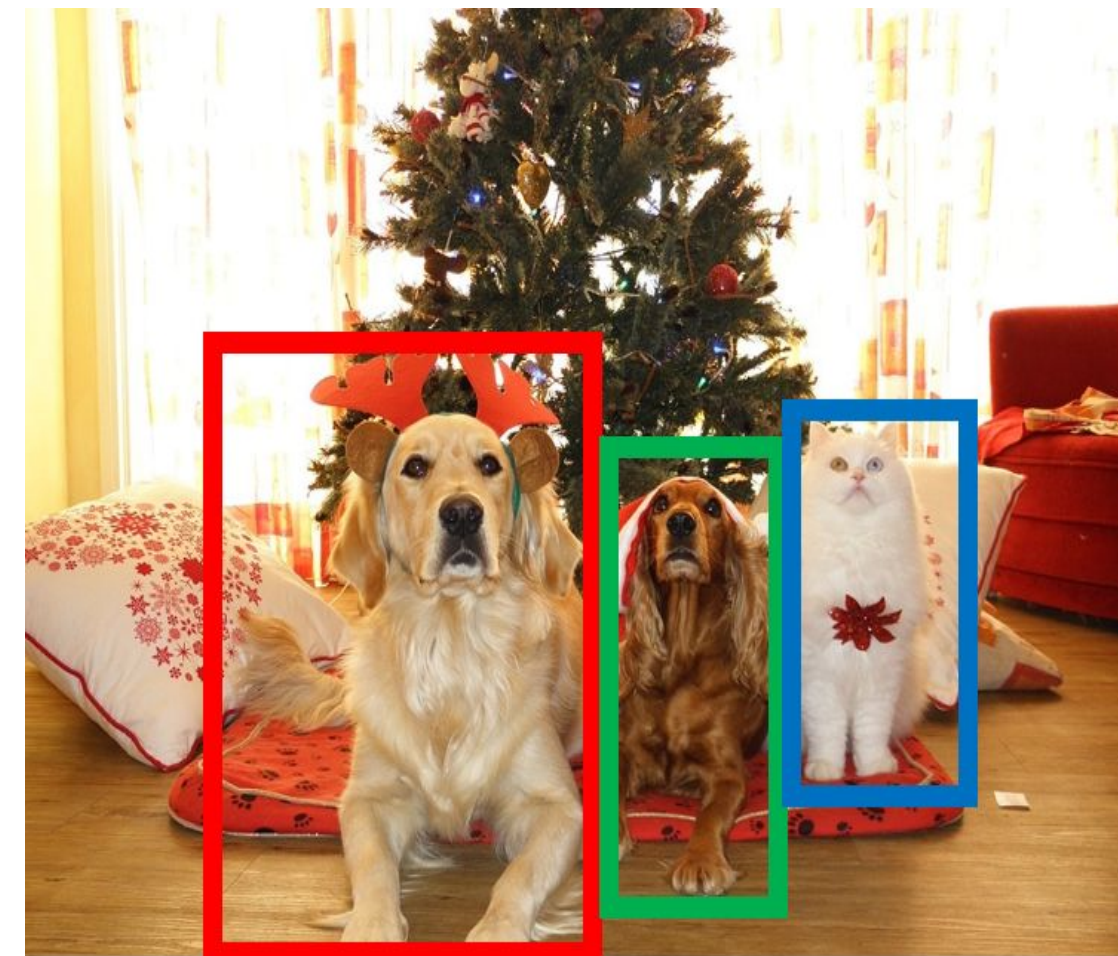**Semantic Segmentation**



**GRASS**, **CAT**, **TREE**, **SKY**

**Object Detection**



**DOG**, **DOG**, **CAT**

**Instance Segmentation**



**DOG**, **DOG**, **CAT**

No spatial extent

No objects, just pixels

Multiple Object

This image is CC0 public domain

Slides from Stanford CS231N: Object Detection and Image Segmentation

But manipulating objects require 3D reasoning!

# Depth cameras give us 3D information!

Left infrared camera

Right infrared camera

Infrared dot matrix projector

RGB camera

Add Source

2D | 3D ⚙

Intel RealSense D435I ⟷ 3.2 ✕
File: "20210622_173905.bag"

⏸ Pause | ↻ Reset | 🔒 Lock | ▣ Source | 🪣 Texture | 🖌 Shading | 📏 Measure | 🐾 Route | 💾 Export

◀ ⏹ ⏸ ▶ ↻  Speed: x1  ⓘ

0:00:03                                    0:00:26

▼ Stereo Module                                       ● on

Resolution:          848 x 480
Frame Rate (FPS):    30
Infrared 2           Y8
Infrared 1           Y8
Depth                Z16

Visual Preset:                              ?
                        0

Emitter Enabled:                            ?
                        1

Enable Auto Exposure:                       ?
                        1

▶ Controls
▶ Depth Visualization
▶ Post-Processing                           ●

▼ RGB Camera                                ● on

Resolution:          1280 x 720
Frame Rate (FPS):    30
Color                RGB8

Enable Auto Exposure:                       ?
                        1

▶ Controls
▶ Post-Processing                           ●

▼ Motion Module                             ● on

Gyro stream:
Format:              MOTION_XYZ32F
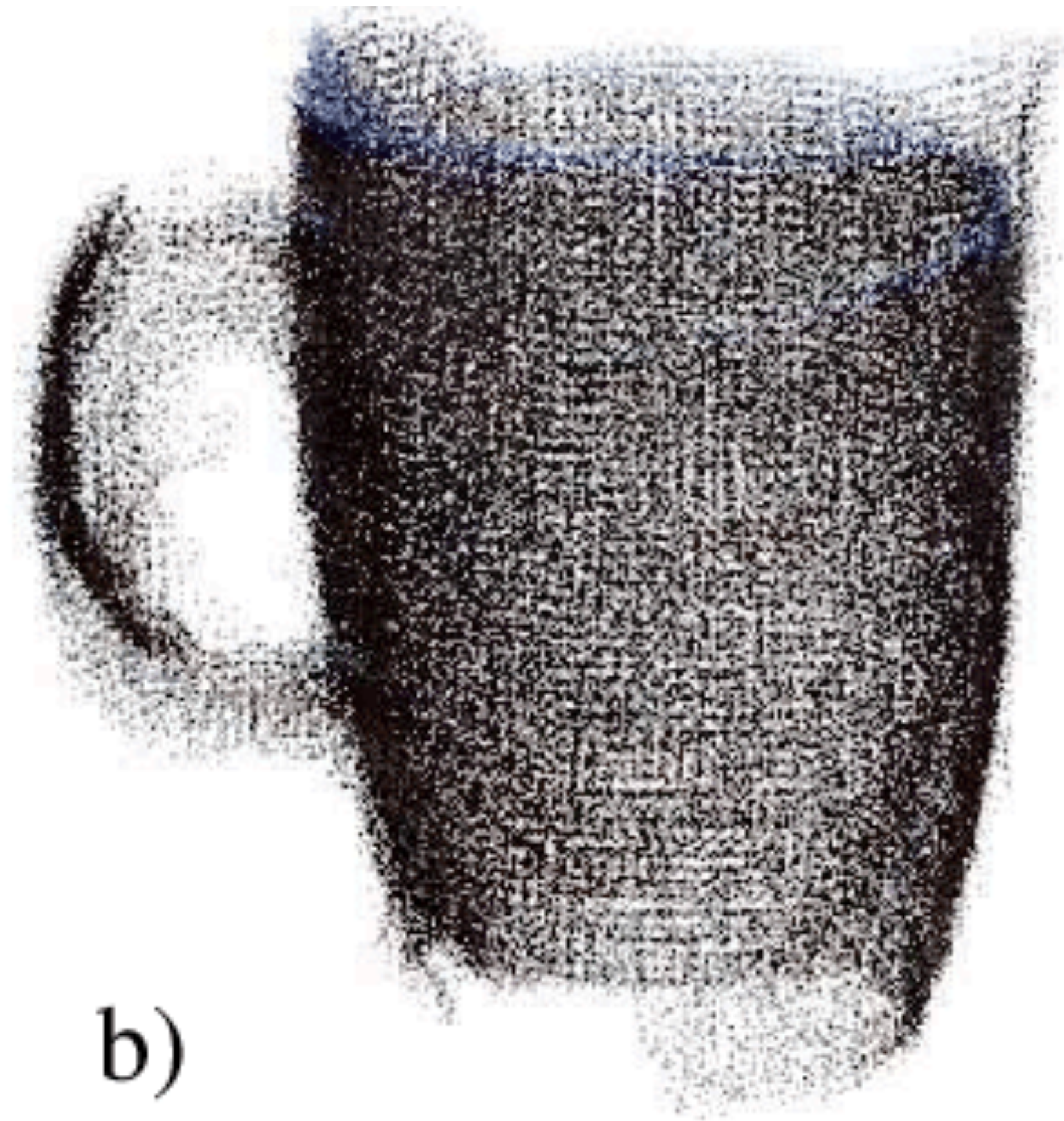Frame Rate (FPS):    200
Accel stream:
Format:              MOTION_XYZ32F
Frame Rate (FPS):    63

▶ Controls

◎2 ⚠ ⓘ80 🔍                                    🔋 ⬛31%

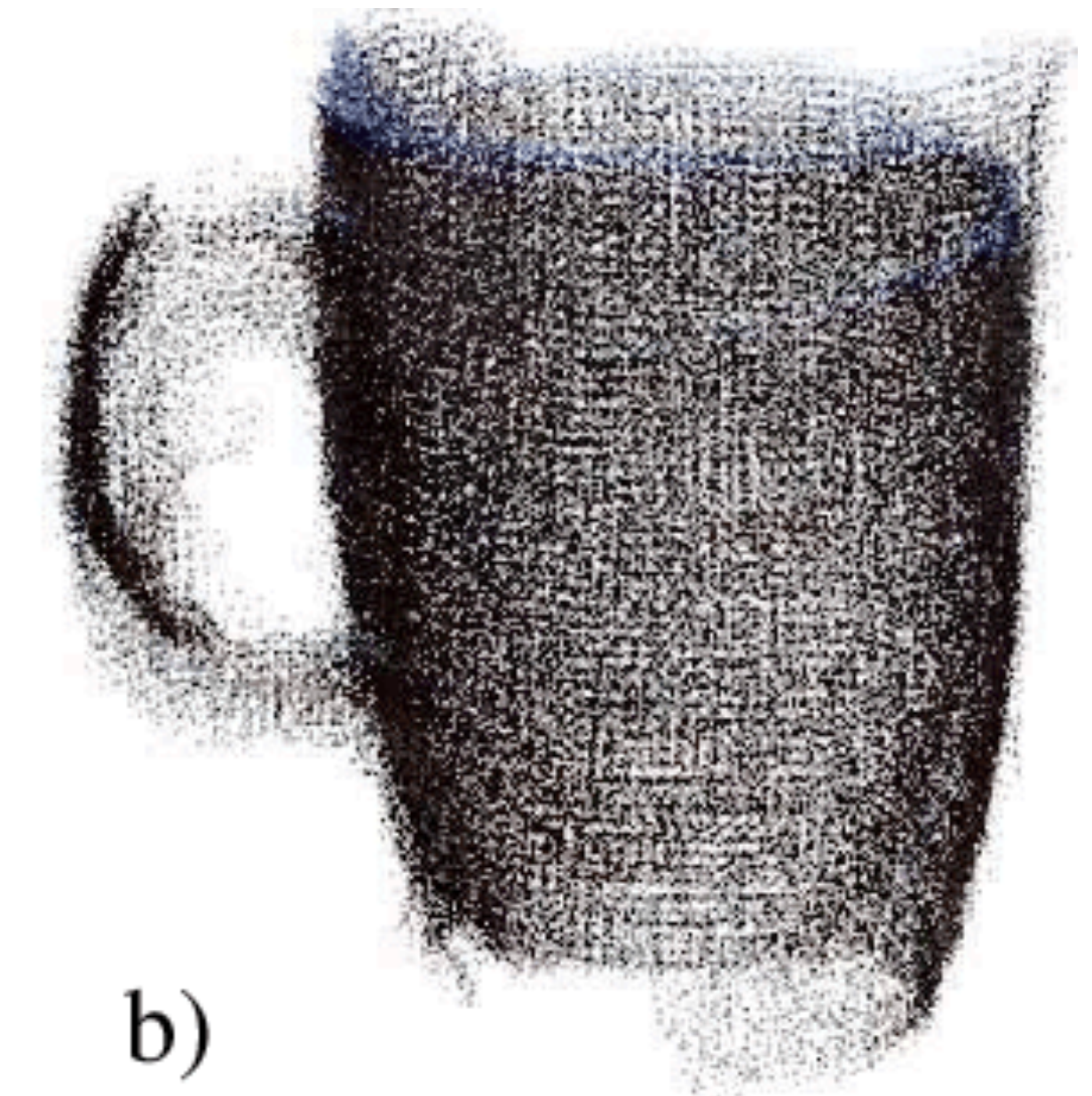# Masked Depth Image -> Point Cloud



a)

b)

# Think-Pair-Share!

Think (30 sec): Given a point cloud of an object, how would you learn where to grasp it? What are some informative features?

Pair: Find a partner

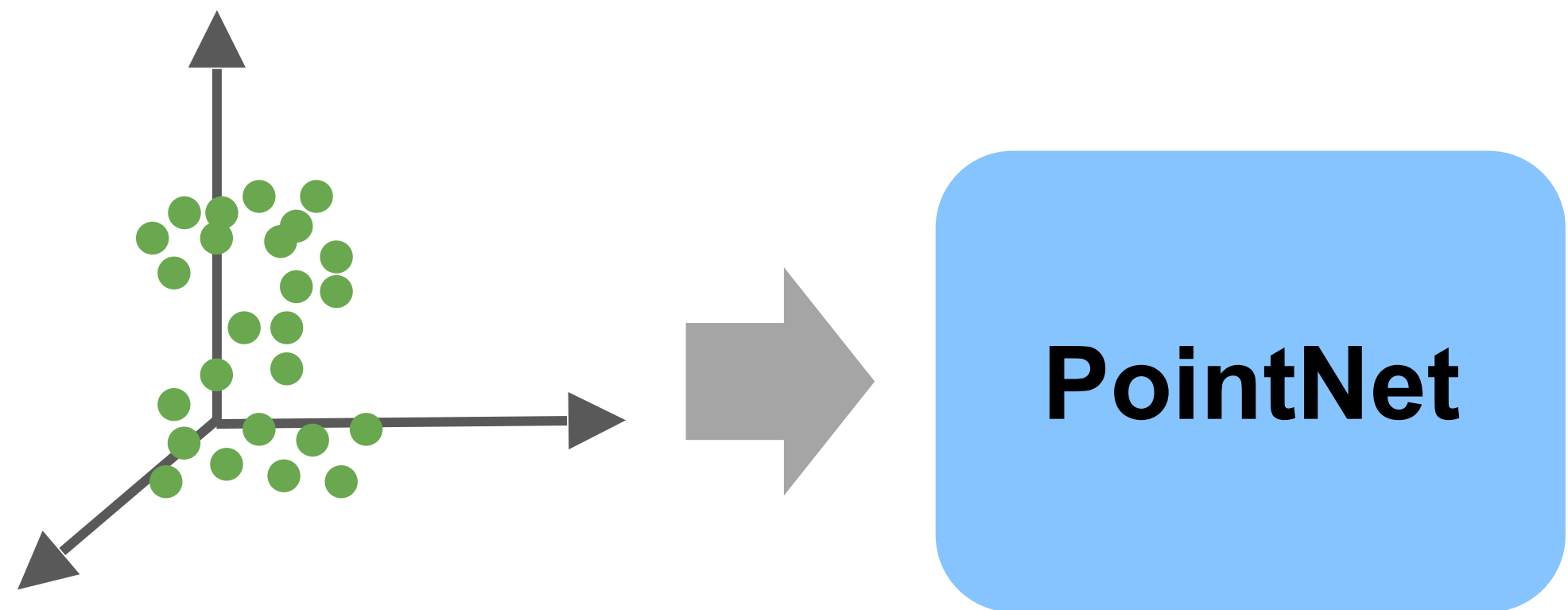Share (45 sec): Partners exchange
                ideas

b)

Can we learn effective feature learning directly on point clouds?
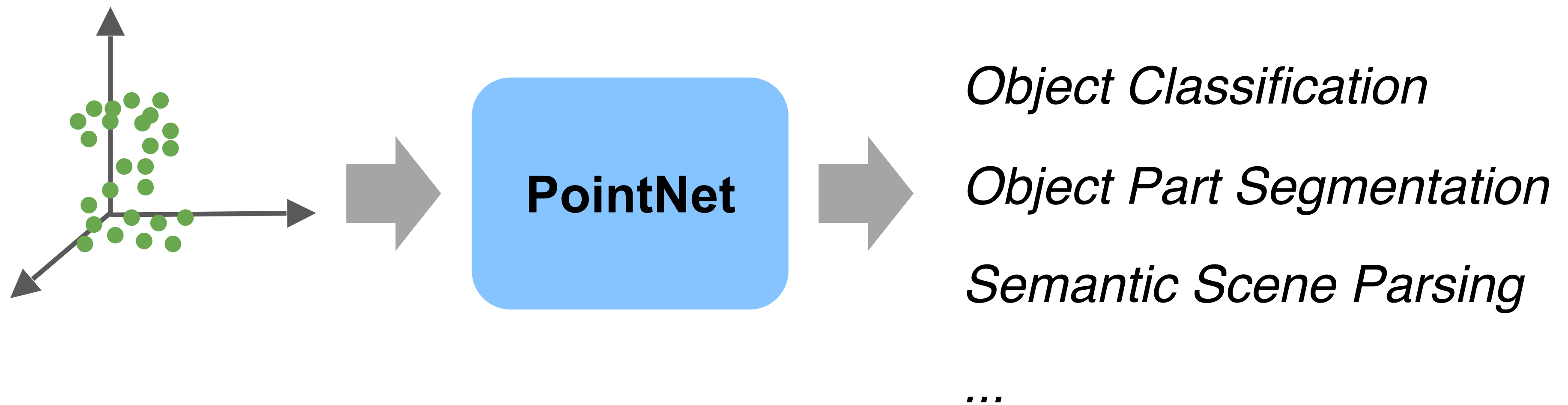
# PointNet

End-to-end learning for **scattered, unordered** point data



**PointNet**

Slides from Qi et al, CVP 2017  http://stanford.edu/~rqi/pointnet/docs/cvpr17_pointnet_slides.pdf
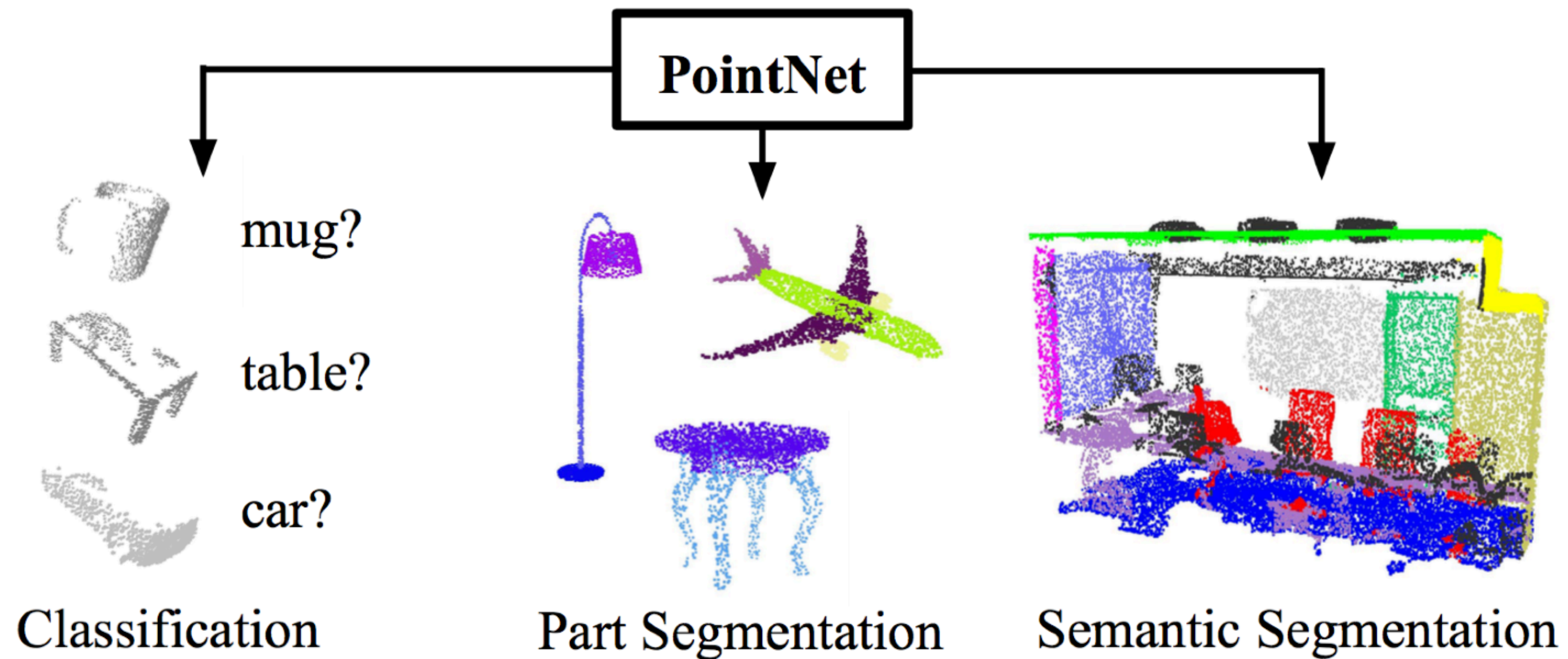
# PointNet

End-to-end learning for **scattered, unordered** point data

**Unified** framework for various tasks



PointNet

*Object Classification*

*Object Part Segmentation*

*Semantic Scene Parsing*

*...*

# PointNet

End-to-end learning for **scattered, unordered** point data

**Unified** framework for various tasks

# Why is learning with point clouds challenging?

# Two Challenges

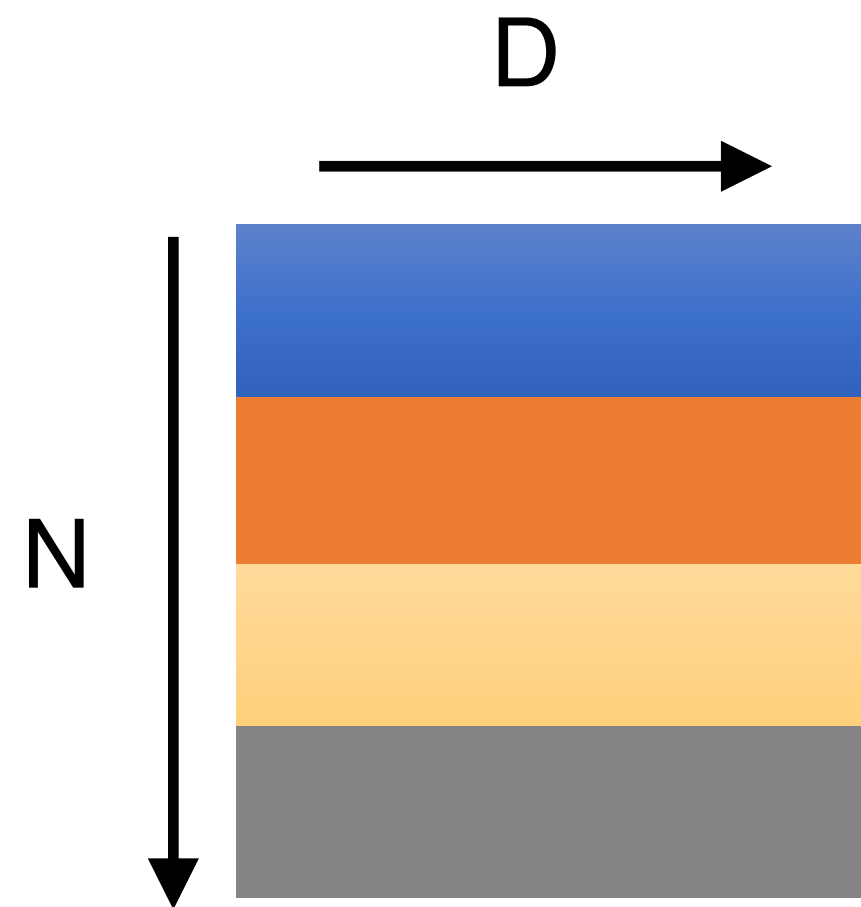Challenge 1: **Unordered** point set as input

*Model needs to be invariant to N! permutation*

Challenge 2: Invariance under **geometric** transformations

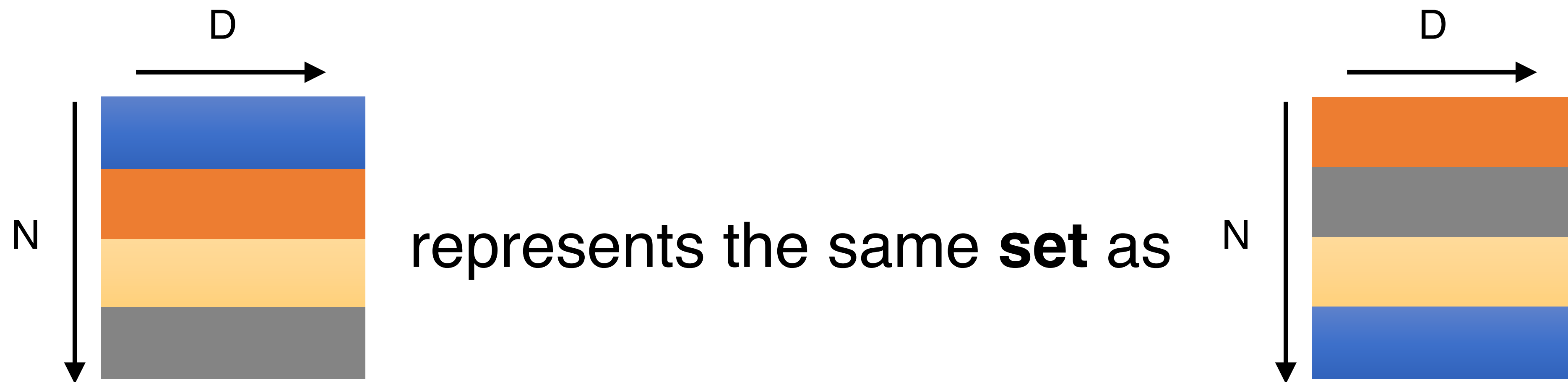*Point cloud rotation should not alter classification results*

# Two Challenges

Challenge 1: **Unordered** point set as input

*Model needs to be invariant to N! permutation*

Challenge 2: Invariance under **geometric** transformations

*Point cloud rotation should not alter classification results*

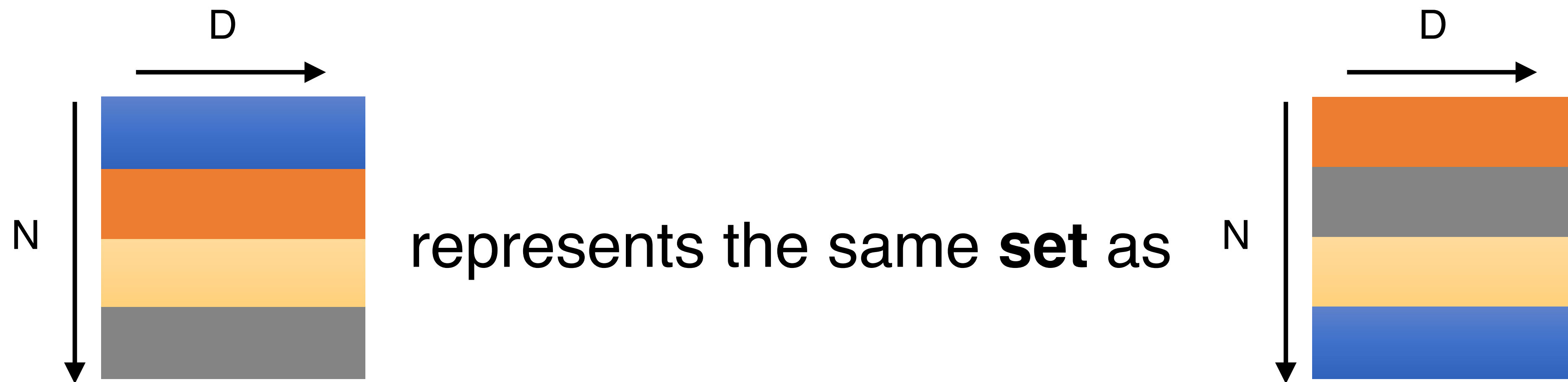Point cloud: N **orderless** points, each represented by a D dim vector

D

N

Point cloud: N **orderless** points, each represented by a D dim vector



represents the same **set** as

Point cloud: N **orderless** points, each represented by a D dim vector



represents the same **set** as

**Model needs to be invariant to N! permutations**

# Permutation Invariance: Symmetric Function

$$f(x_1, x_2, \ldots, x_n) \equiv f(x_{\pi_1}, x_{\pi_2}, \ldots, x_{\pi_n}), \quad x_i \in \mathbb{R}^D$$

# Permutation Invariance: Symmetric Function

$$f(x_1, x_2, \ldots, x_n) \equiv f(x_{\pi_1}, x_{\pi_2}, \ldots, x_{\pi_n}), \quad x_i \in \mathbb{R}^D$$

**Examples:**

$$f(x_1, x_2, \ldots, x_n) = \max\{x_1, x_2, \ldots, x_n\}$$

$$f(x_1, x_2, \ldots, x_n) = x_1 + x_2 + \ldots + x_n$$

…

Slides from Qi et al, CVP 2017  http://stanford.edu/~rqi/pointnet/docs/cvpr17_pointnet_slides.pdf

$$f(x_1, x_2, \ldots, x_n) \equiv f(x_{\pi_1}, x_{\pi_2}, \ldots, x_{\pi_n}), \quad x_i \in \mathbb{R}^D$$

## Examples:

$$f(x_1, x_2, \ldots, x_n) = \max\{x_1, x_2, \ldots, x_n\}$$

$$f(x_1, x_2, \ldots, x_n) = x_1 + x_2 + \ldots + x_n$$

…

**How can we construct a family of symmetric functions by neural networks?**

Slides from Qi et al, CVP 2017  http://stanford.edu/~rqi/pointnet/docs/cvpr17_pointnet_slides.pdf
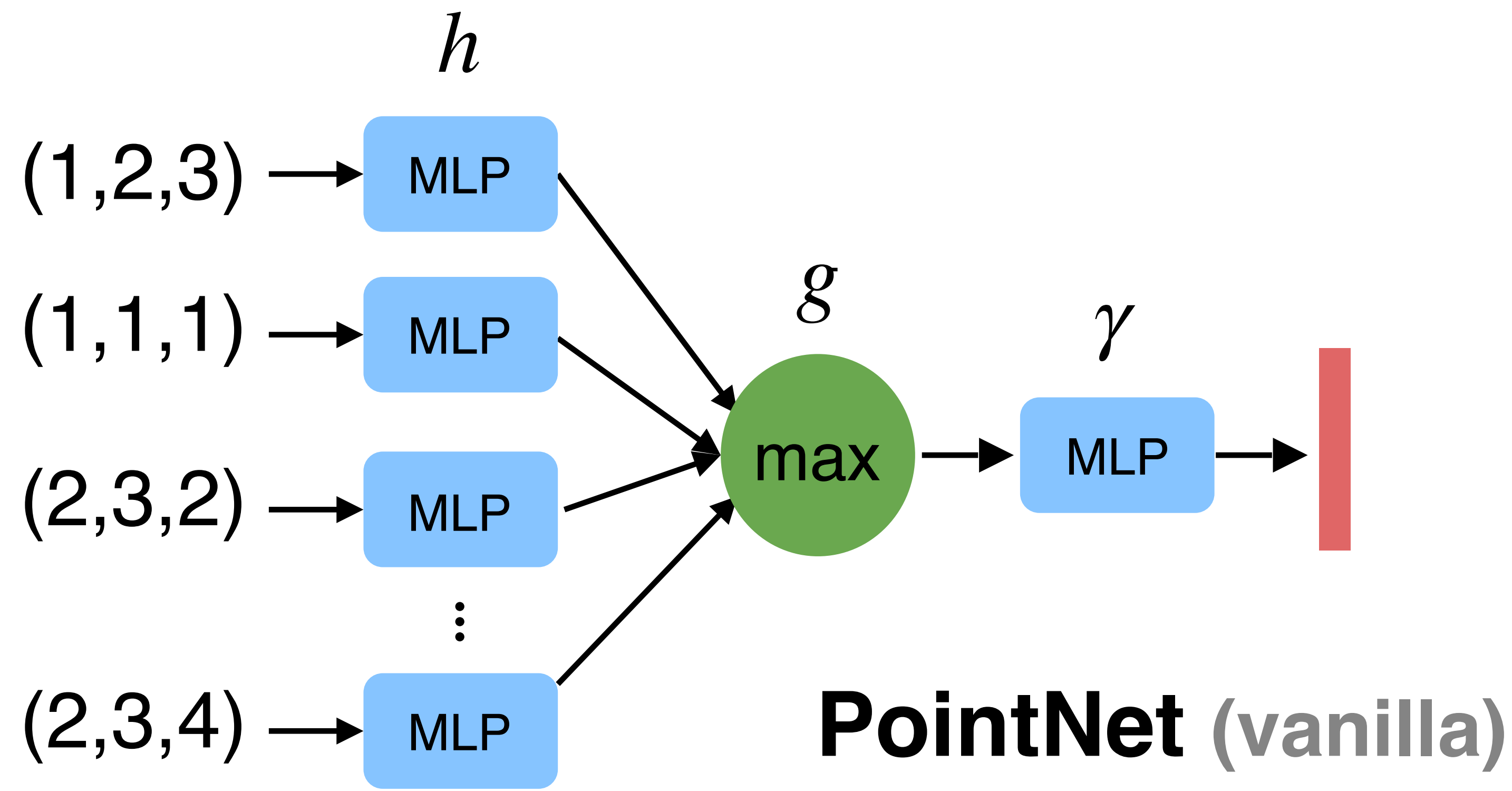
# Permutation Invariance: Symmetric Function

**Observe:**

$f(x_1, x_2, \ldots, x_n) = \gamma \circ g(h(x_1), \ldots, h(x_n))$ is symmetric if $g$ is symmetric

**Observe:**

$$f(x_1, x_2, \ldots, x_n) = \gamma \circ g(h(x_1), \ldots, h(x_n))$$ is symmetric if $g$ is symmetric

$h$

$(1,2,3) \longrightarrow$ 

$(1,1,1) \longrightarrow$ 

$(2,3,2) \longrightarrow$ 

$\vdots$

$(2,3,4) \longrightarrow$ 

Slides from Qi et al, CVP 2017   http://stanford.edu/~rqi/pointnet/docs/cvpr17_pointnet_slides.pdf

**Observe:**

$$f(x_1, x_2, \ldots, x_n) = \gamma \circ g(h(x_1), \ldots, h(x_n)) \text{ is symmetric if } g \text{ is symmetric}$$



$h$

(1,2,3) →

(1,1,1) →

(2,3,2) →

⋮

(2,3,4) →

$g$

simple symmetric function

Slides from Qi et al, CVP 2017   http://stanford.edu/~rqi/pointnet/docs/cvpr17_pointnet_slides.pdf

**Observe:**

$$f(x_1, x_2, \ldots, x_n) = \gamma \circ g(h(x_1), \ldots, h(x_n)) \text{ is symmetric if } g \text{ is symmetric}$$



$h$

(1,2,3) →

(1,1,1) →

(2,3,2) →

⋮

(2,3,4) →

simple symmetric function

$g$    $\gamma$

**PointNet** (vanilla)

Slides from Qi et al, CVP 2017   http://stanford.edu/~rqi/pointnet/docs/cvpr17_pointnet_slides.pdf

Empirically, we use **multi-layer perceptron (MLP)** and **max pooling**:



$h$

$(1,2,3) \longrightarrow$ MLP

$(1,1,1) \longrightarrow$ MLP

$(2,3,2) \longrightarrow$ MLP

$\vdots$

$(2,3,4) \longrightarrow$ MLP

$g$

max

$\gamma$

MLP

**PointNet** (vanilla)

Slides from Qi et al, CVP 2017   http://stanford.edu/~rqi/pointnet/docs/cvpr17_pointnet_slides.pdf

# Two Challenges

Challenge 1: **Unordered** point set as input

*Model needs to be invariant to N! permutation*

## Challenge 2: Invariance under **geometric** transformations

*Point cloud rotation should not alter classification results*

# Input Alignment by Transformer Network

Idea: Data dependent transformation for automatic alignment

The transformation is just matrix multiplication!

Slides from Qi et al, CVP 2017  http://stanford.edu/~rqi/pointnet/docs/cvpr17_pointnet_slides.pdf

# Results on Semantic Scene Parsing



Input

Output

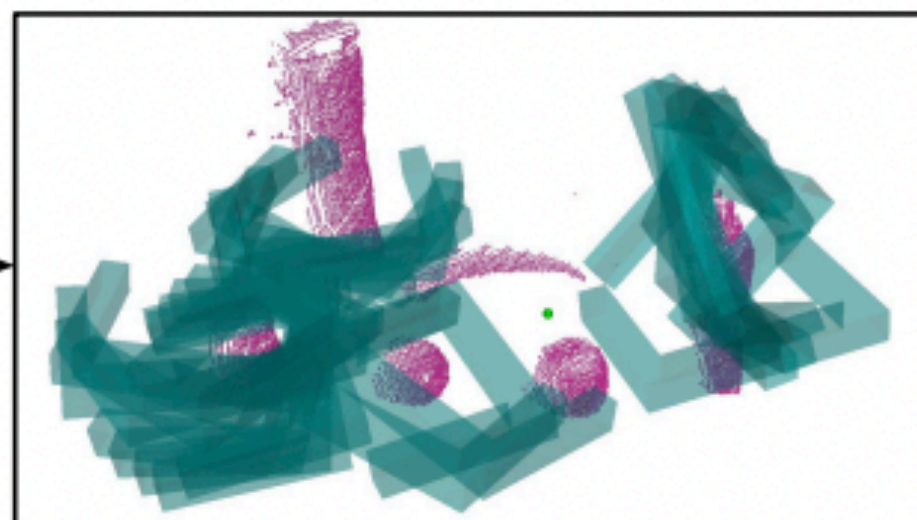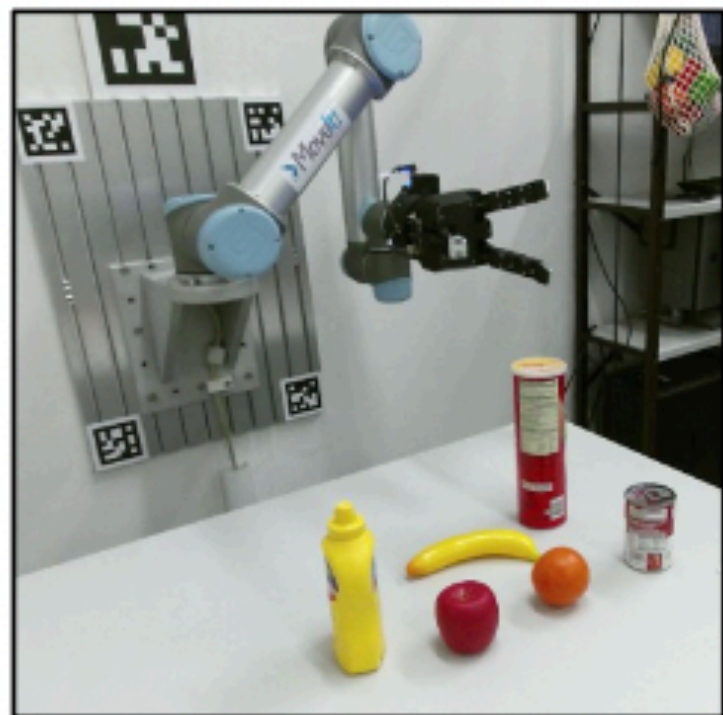Slides from Qi et al, CVP 2017   http://stanford.edu/~rqi/pointnet/docs/cvpr17_pointnet_slides.pdf

How do we use this for learning grasping?

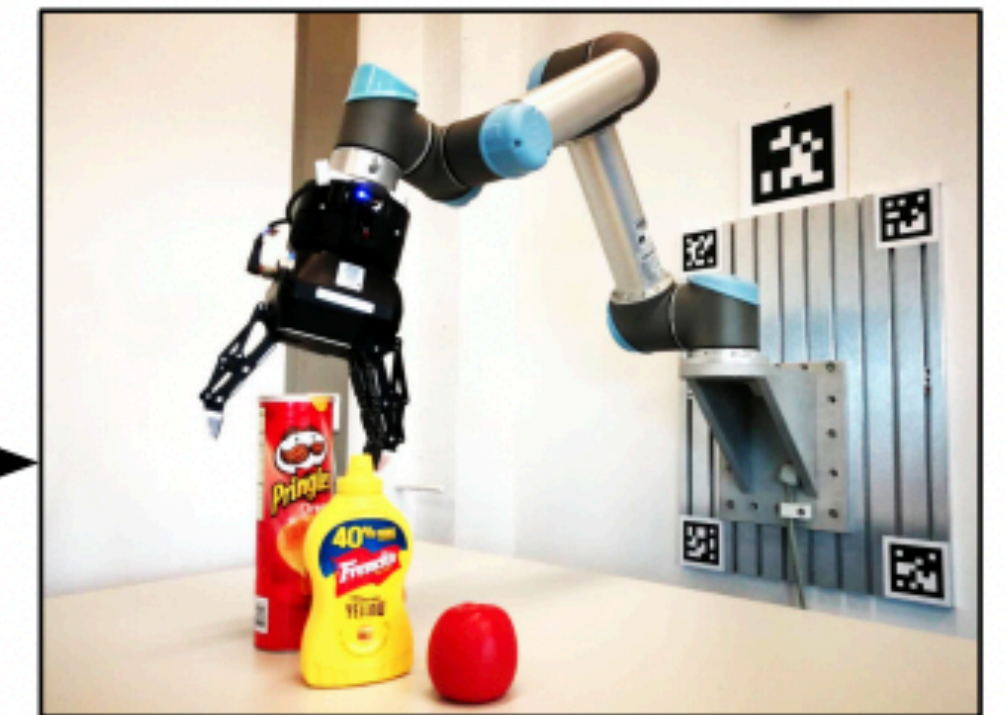# PointNetGPD: Detecting Grasp Configurations from Point Sets



**Grasp Dataset**

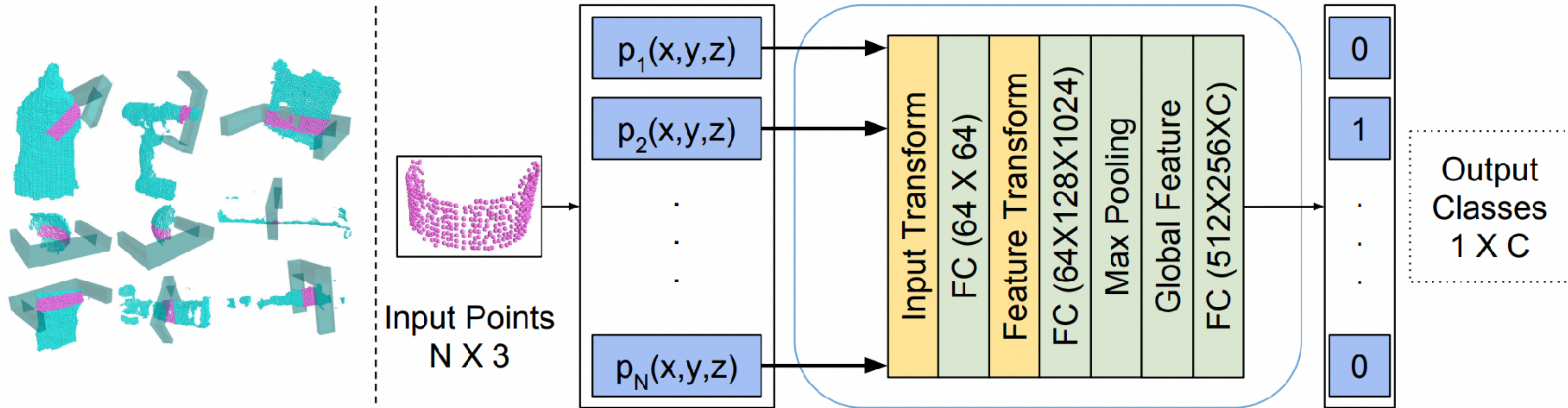**Robot Initial State** **Grasp Candidates Generation**

**Quality Evaluation with PointNet**

**Best Grasp**

**Executed Grasp**

Liang et al.

# PointNetGPD: Detecting Grasp Configurations from Point Sets



Input Points N X 3

$p_1(x,y,z)$

$p_2(x,y,z)$

$p_N(x,y,z)$

Input Transform

FC (64 X 64)

Feature Transform

FC (64X128X1024)

Max Pooling

Global Feature

FC (512X256XC)

0
1
0

Output Classes 1 X C

Liang et al.

But ... what if we don't have RGBD data?

# Reasons for not having depth data

Don't have a depth sensor!

Outside in the sun / beyond maximum range

Glass or transparent objects

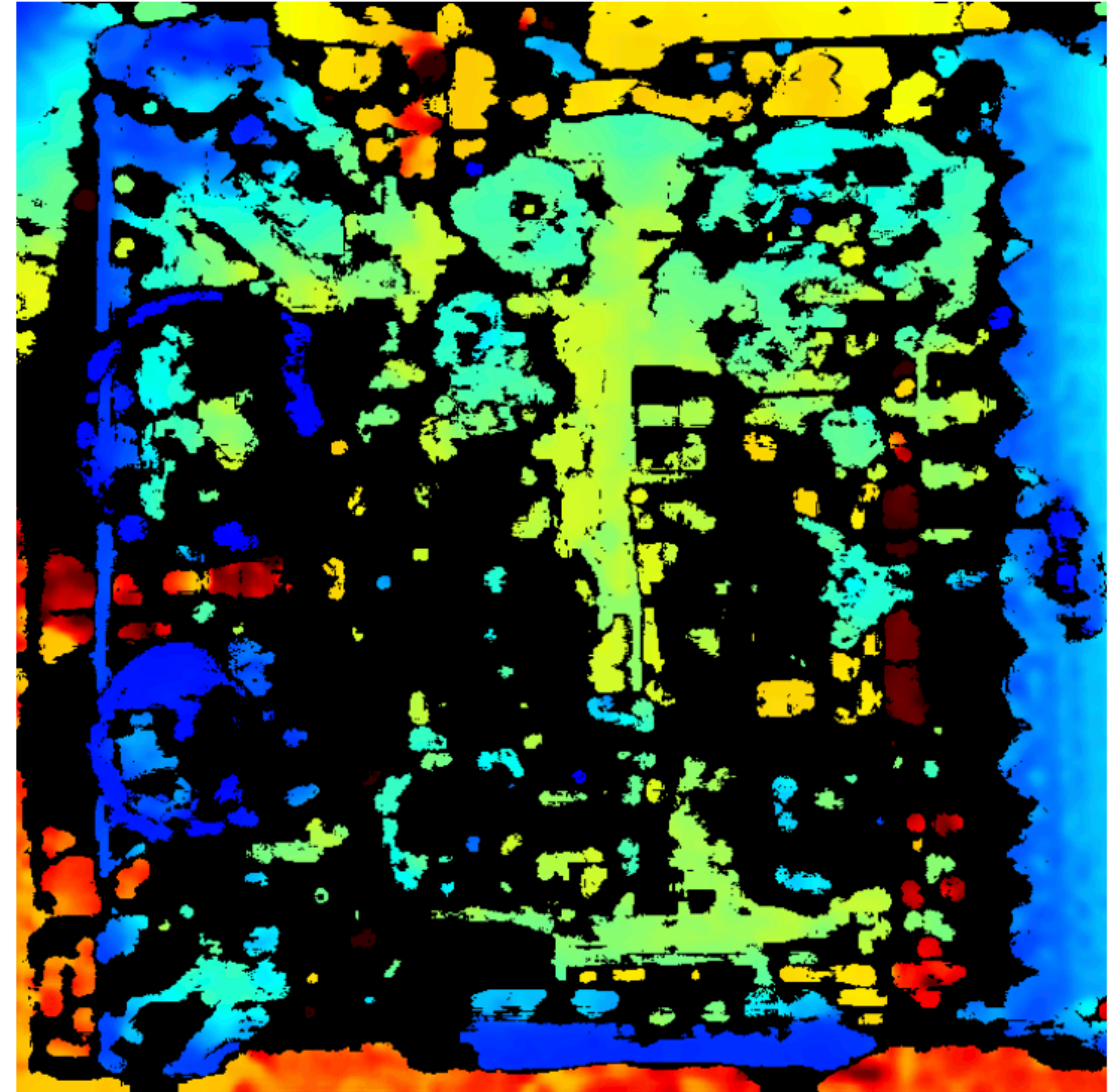# Glass or transparent objects

**Real-world Scene**

**RealSense D410 Depth Image**
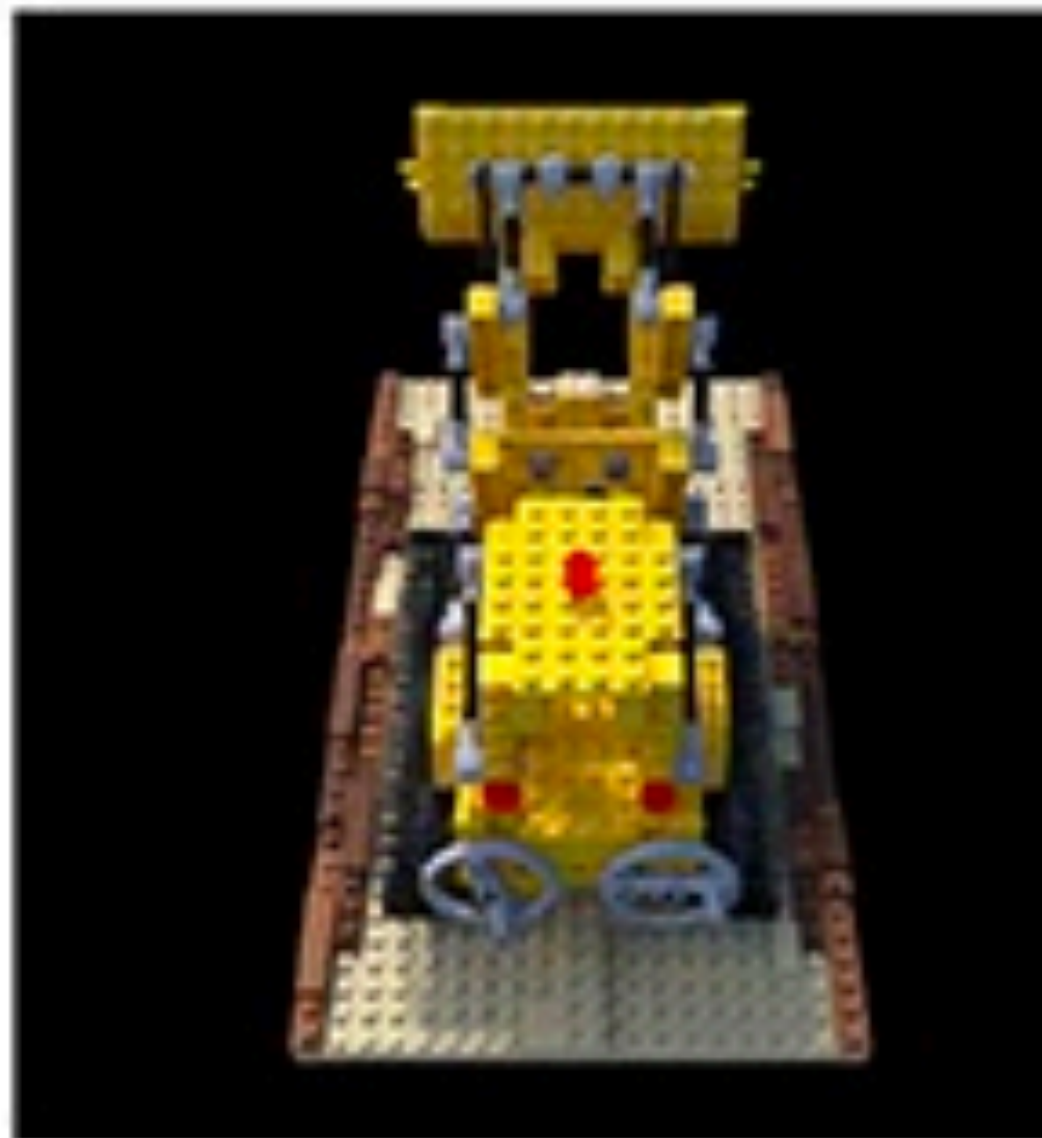
# Glass or transparent objects



**Dishwasher Real-world Scene**

**RealSense D410 Depth Image** 🔗

43

# Let's say I just have a set of images & camera poses



$$x_1, y_1, z_1, \theta_1, \phi_1 \qquad x_2, y_2, z_2, \theta_2, \phi_2 \qquad x_3, y_3, z_3, \theta_3, \phi_3$$
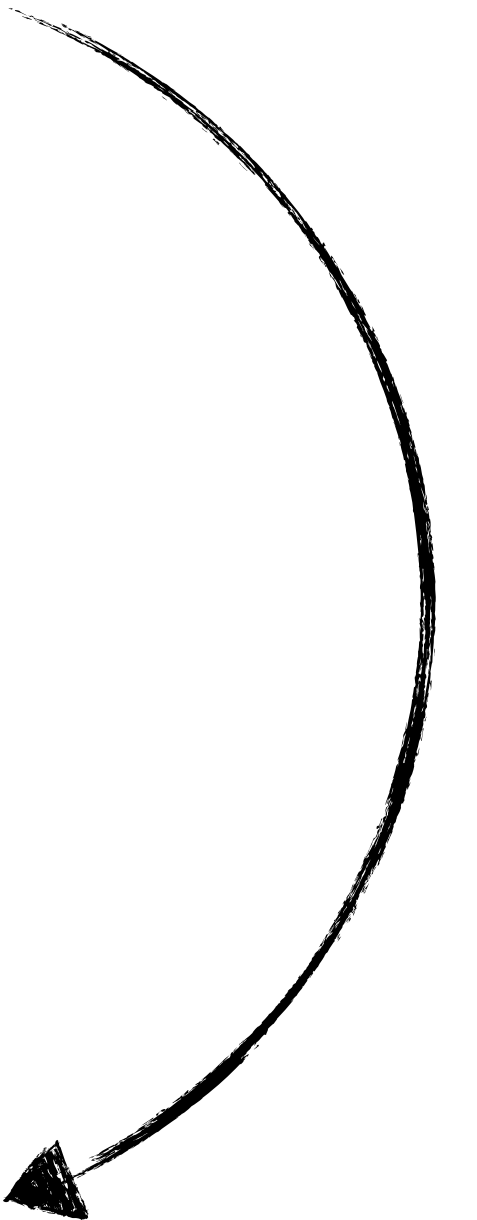
# How do we predict a 3D structure?

2D images



Camera Poses

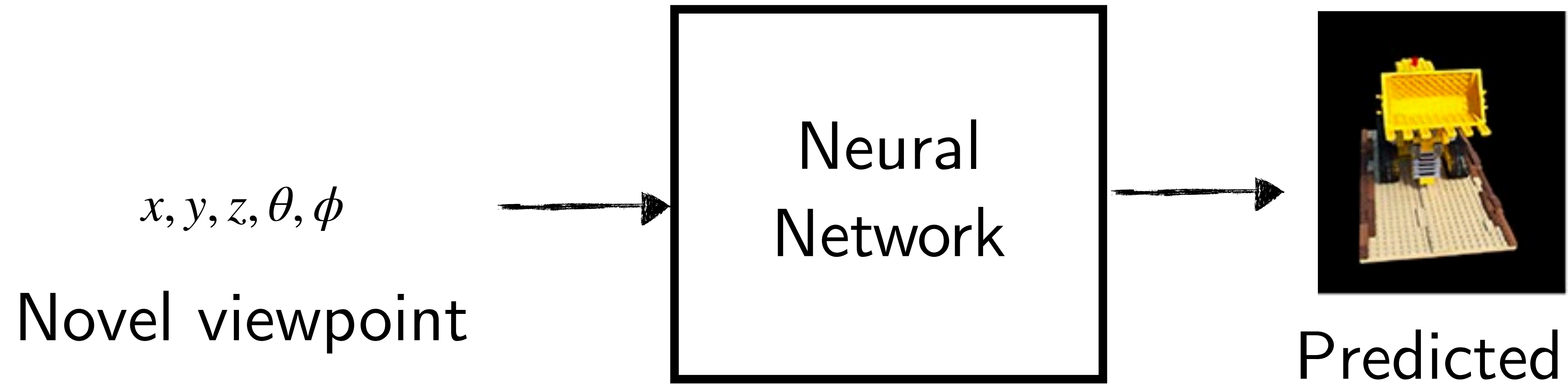$x_1, y_1, z_1, \theta_1, \phi_1$ $\qquad$ $x_2, y_2, z_2, \theta_2, \phi_2$ $\qquad$ $x_3, y_3, z_3, \theta_3, \phi_3$
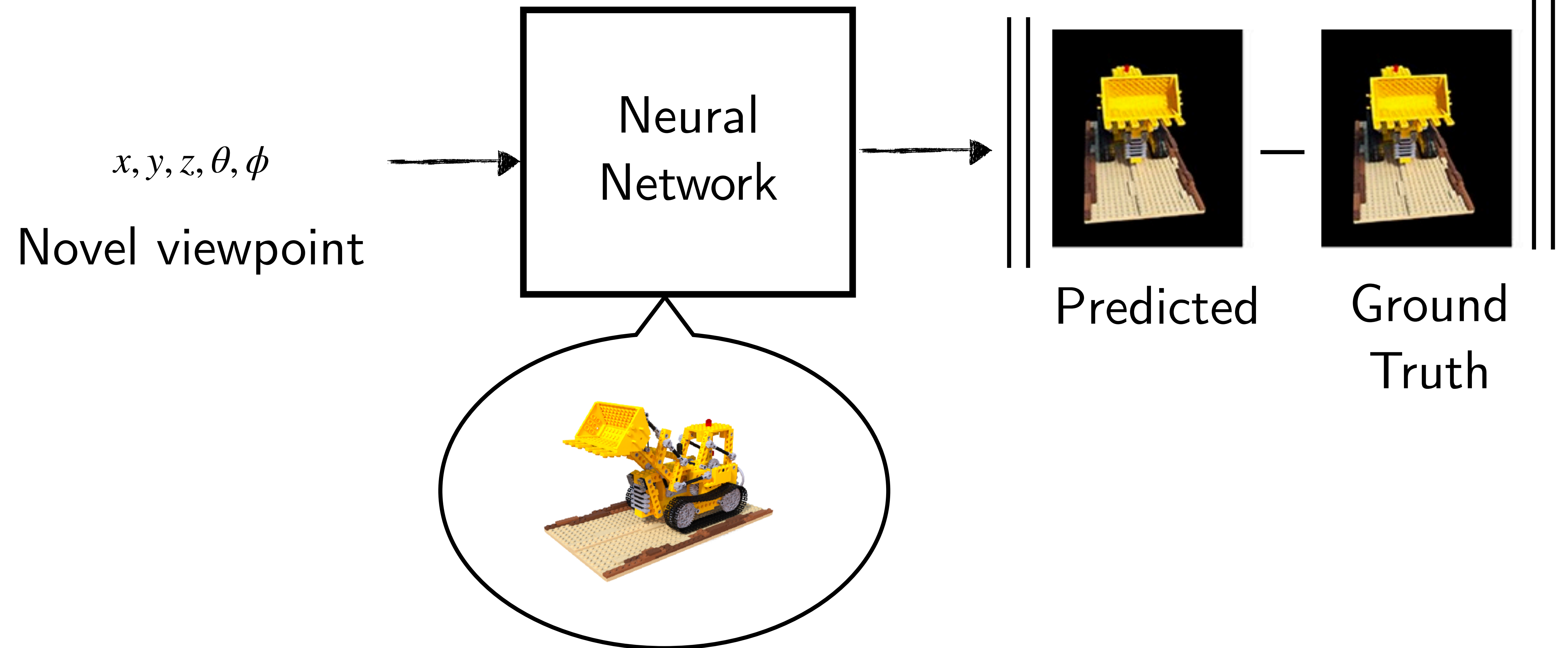
3D structure

If you can predict how the object will look from a *novel viewpoint,* you have *implicitly* modeled the 3D structure

# Let's setup a learning problem

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

**Novel viewpoint**

Neural Network

Predicted

# Let's setup a learning problem



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

Novel viewpoint

Neural Network

Predicted    Ground Truth

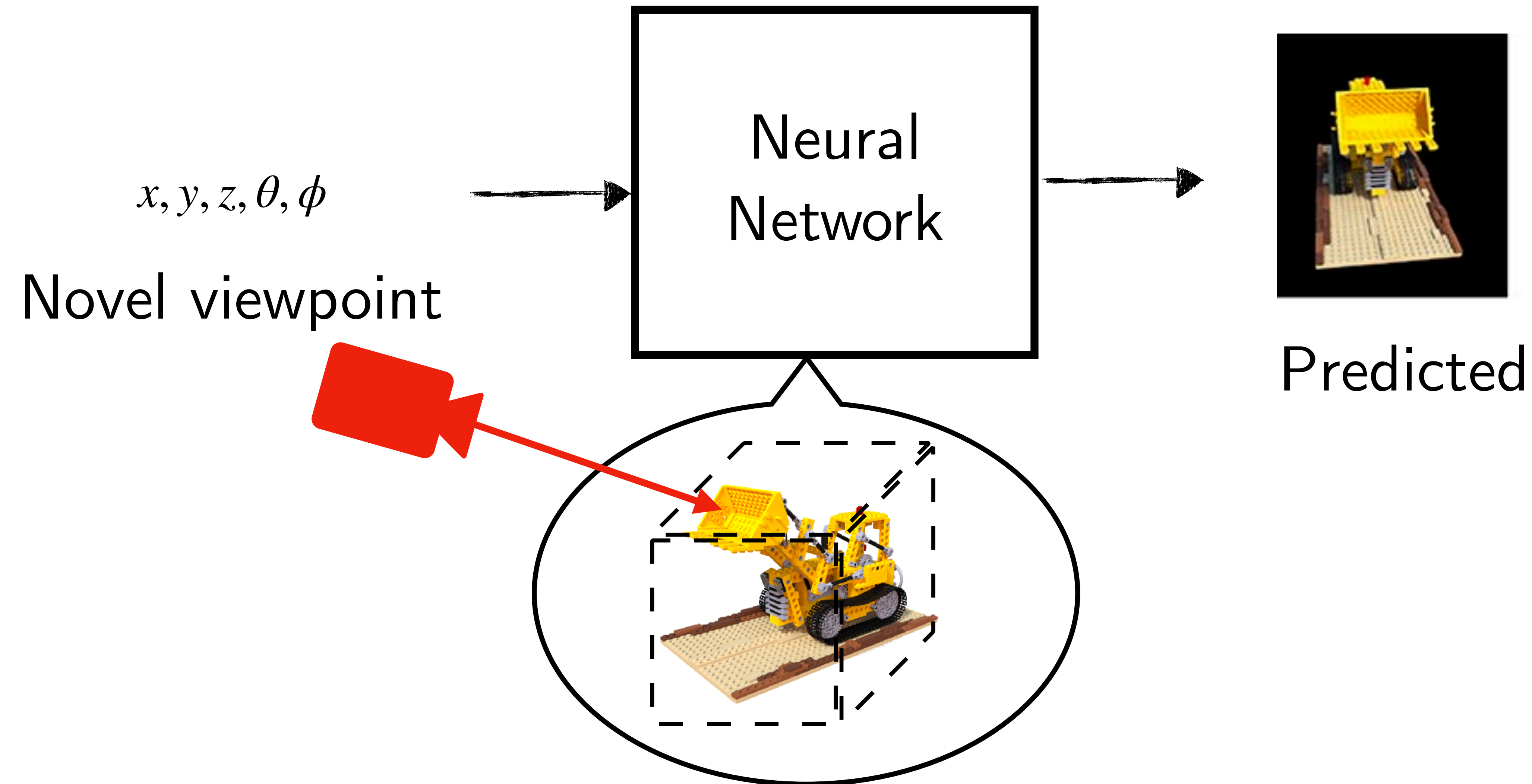# Simple idea:

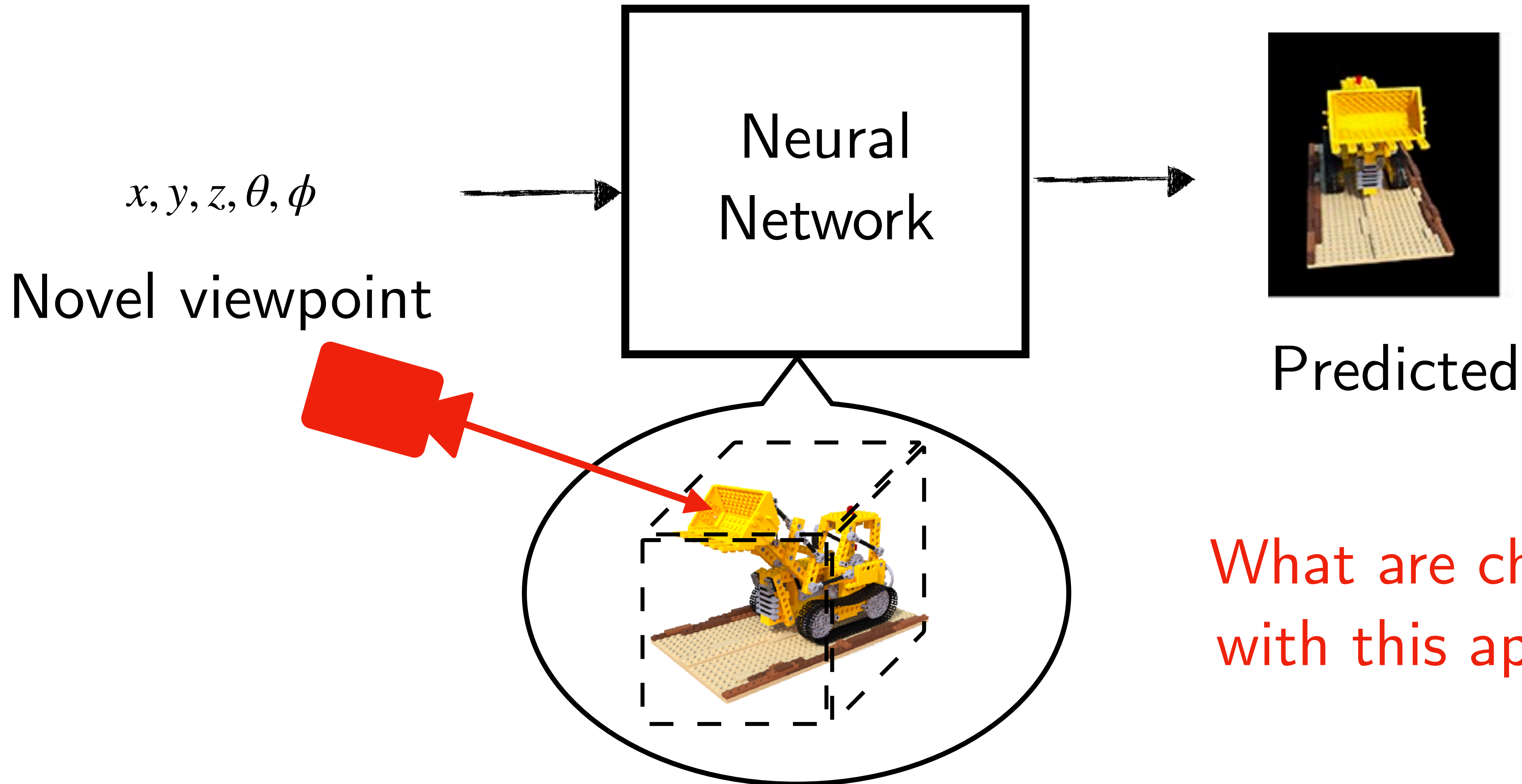## Can't we just make the neural network predict a 3D voxel grid of RGB values?

# Predict a 3D grid, render image from viewpoint

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

Novel viewpoint

Neural Network

Predicted

# Predict a 3D grid, render image from viewpoint

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

Novel viewpoint

Neural Network

Predicted

What are challenges with this approach?

# Challenges

Discretization loses information!
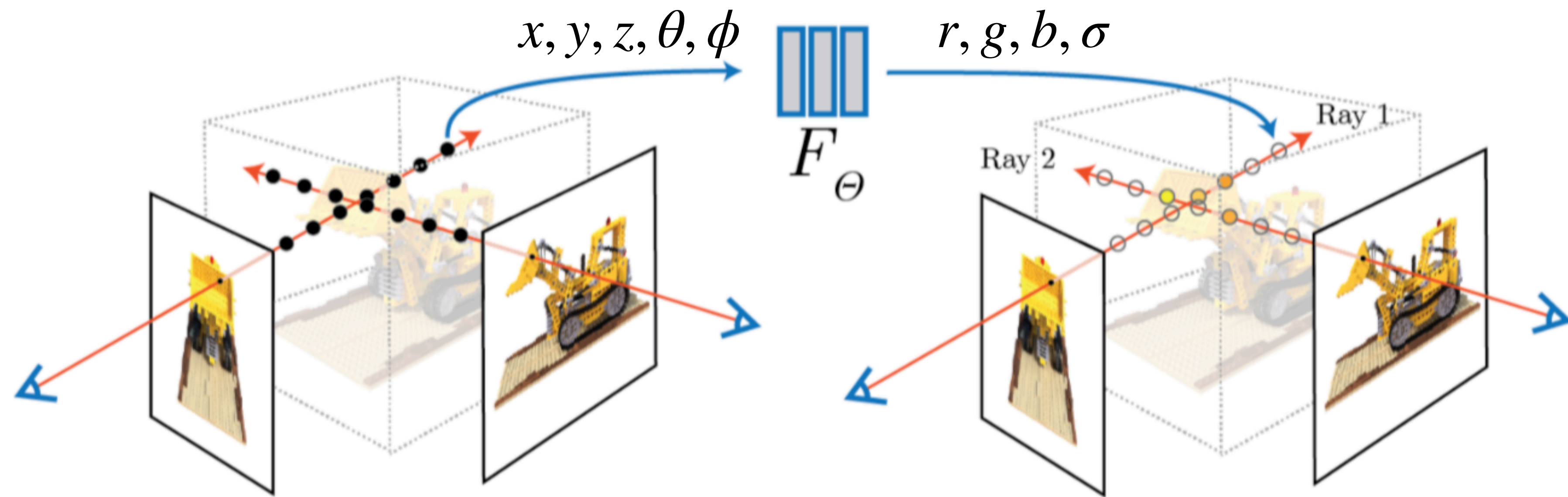
Memory Inefficient!

Not differentiable! (Not a continuous projection)

NERF to the rescue!

# What are Neural Radiance Fields (NeRFs)?

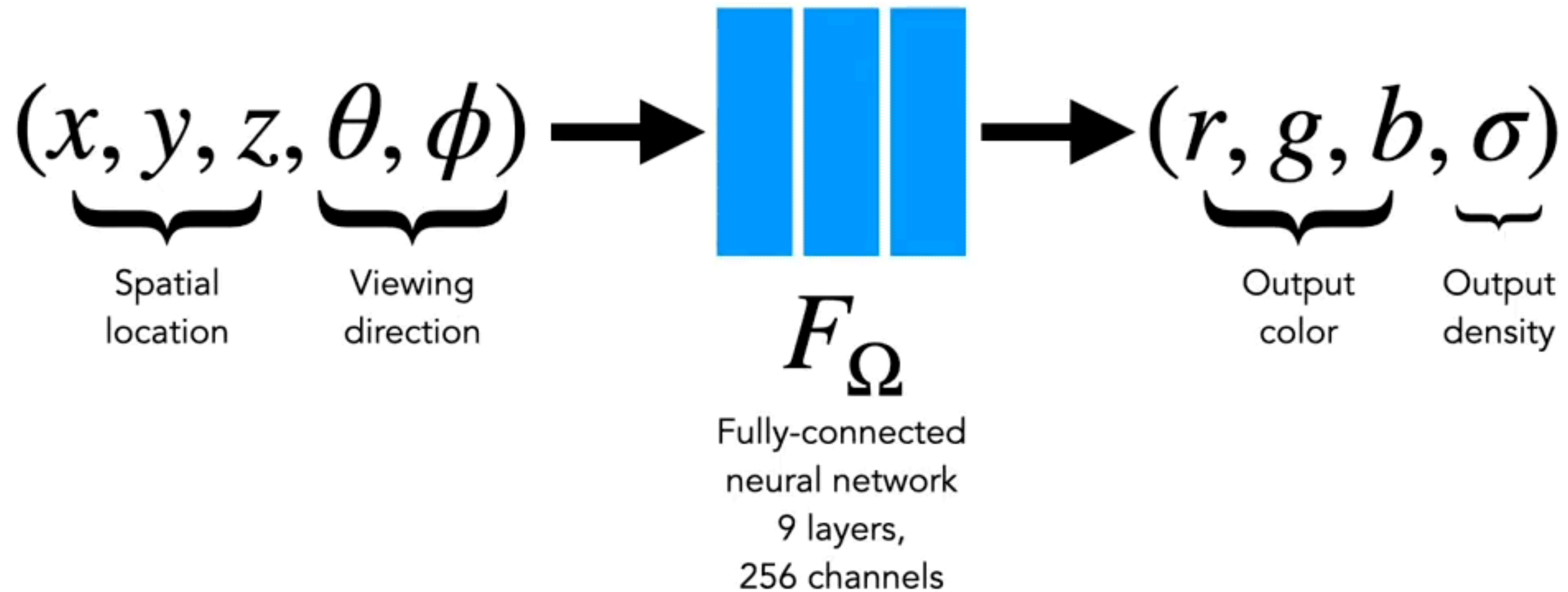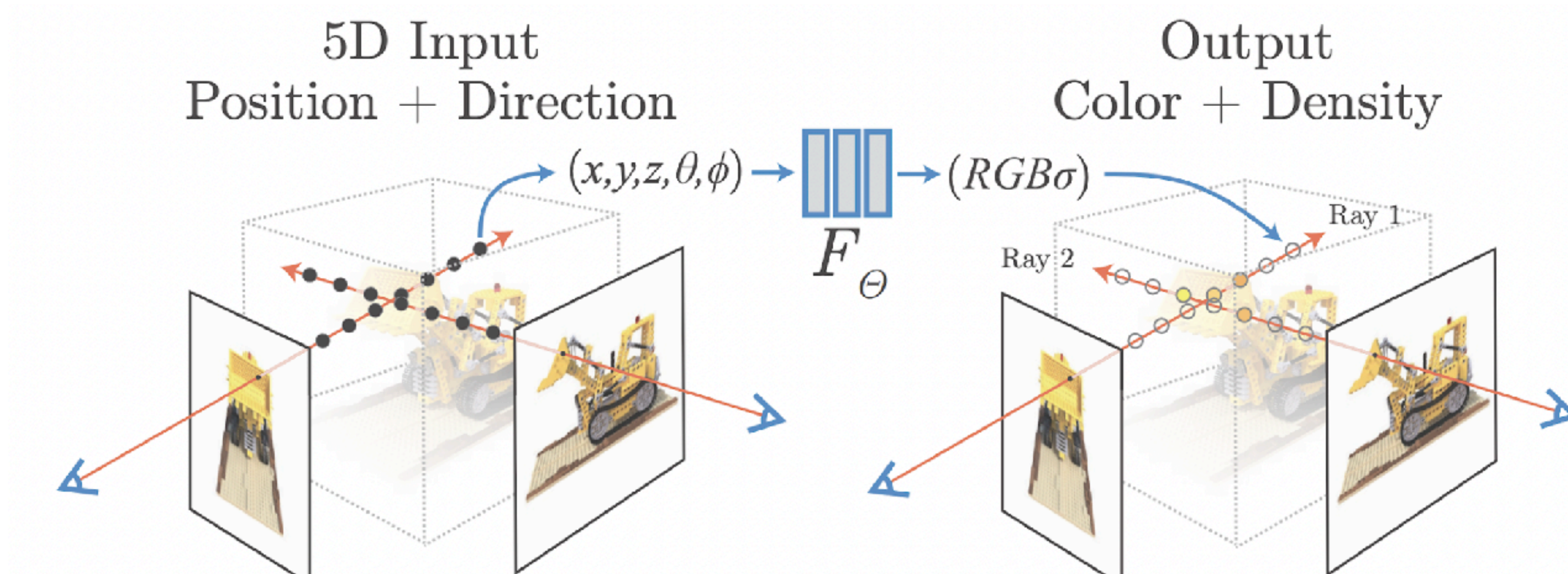Idea: Use a neural network to *implicitly* represent 3D volume!



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

$F_{\Theta}$

$r, g, b, \sigma$

✓ No Discretization    ✓ Compressible    ✓ Differentiable

# NeRF scene representation

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

Spatial location | Viewing direction

$F_\Omega$

Fully-connected neural network 9 layers, 256 channels

Output color | Output density

# Differentiable Loss Function



$$\min_{\Theta} \sum_i ||\mathbf{render}_i(F_{\Theta}) - I_i^{gt}||^2$$
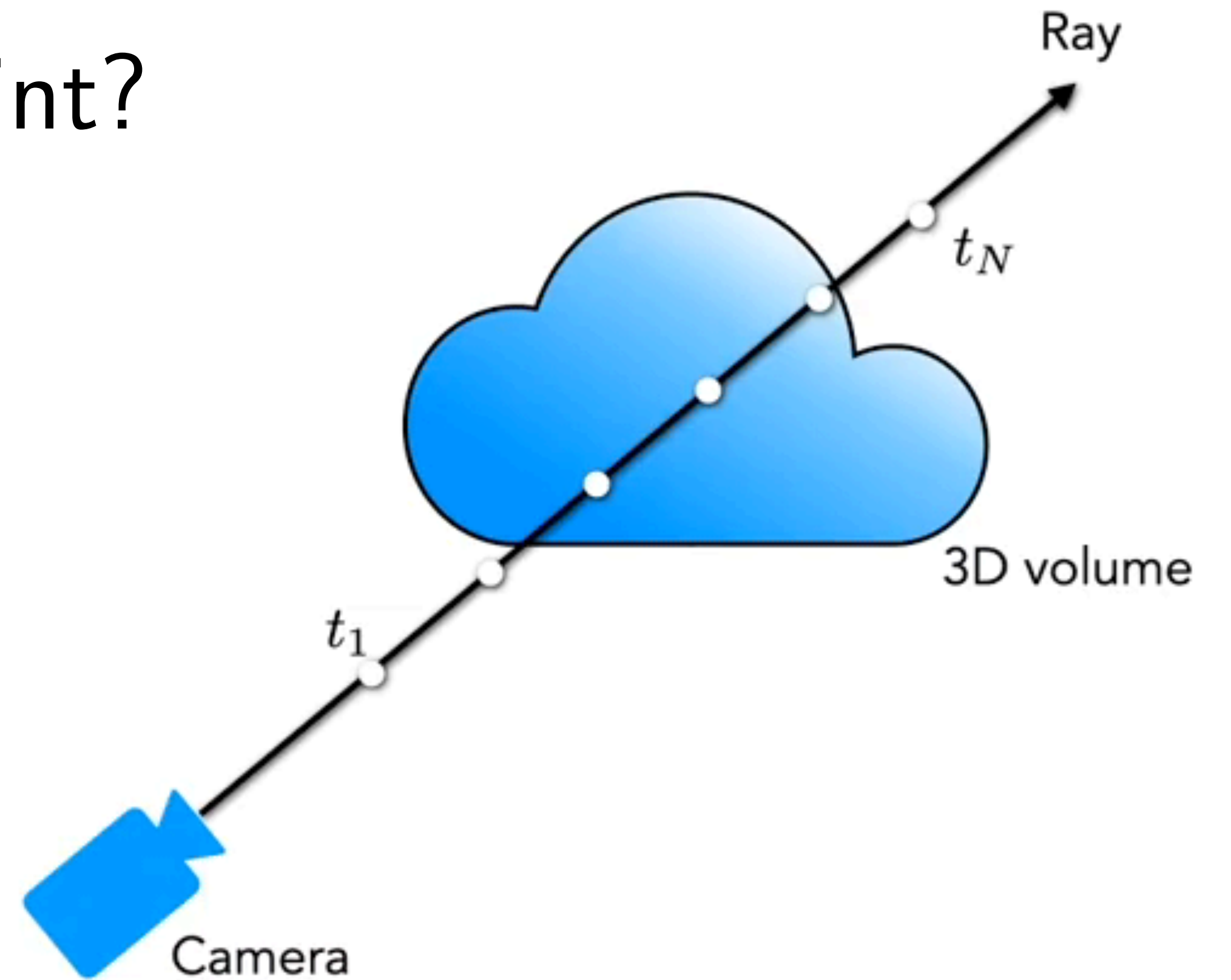
# What is the render() function?
# How is it differentiable?

# Volume rendering model

What is the color from this viewpoint?

# Volume rendering model

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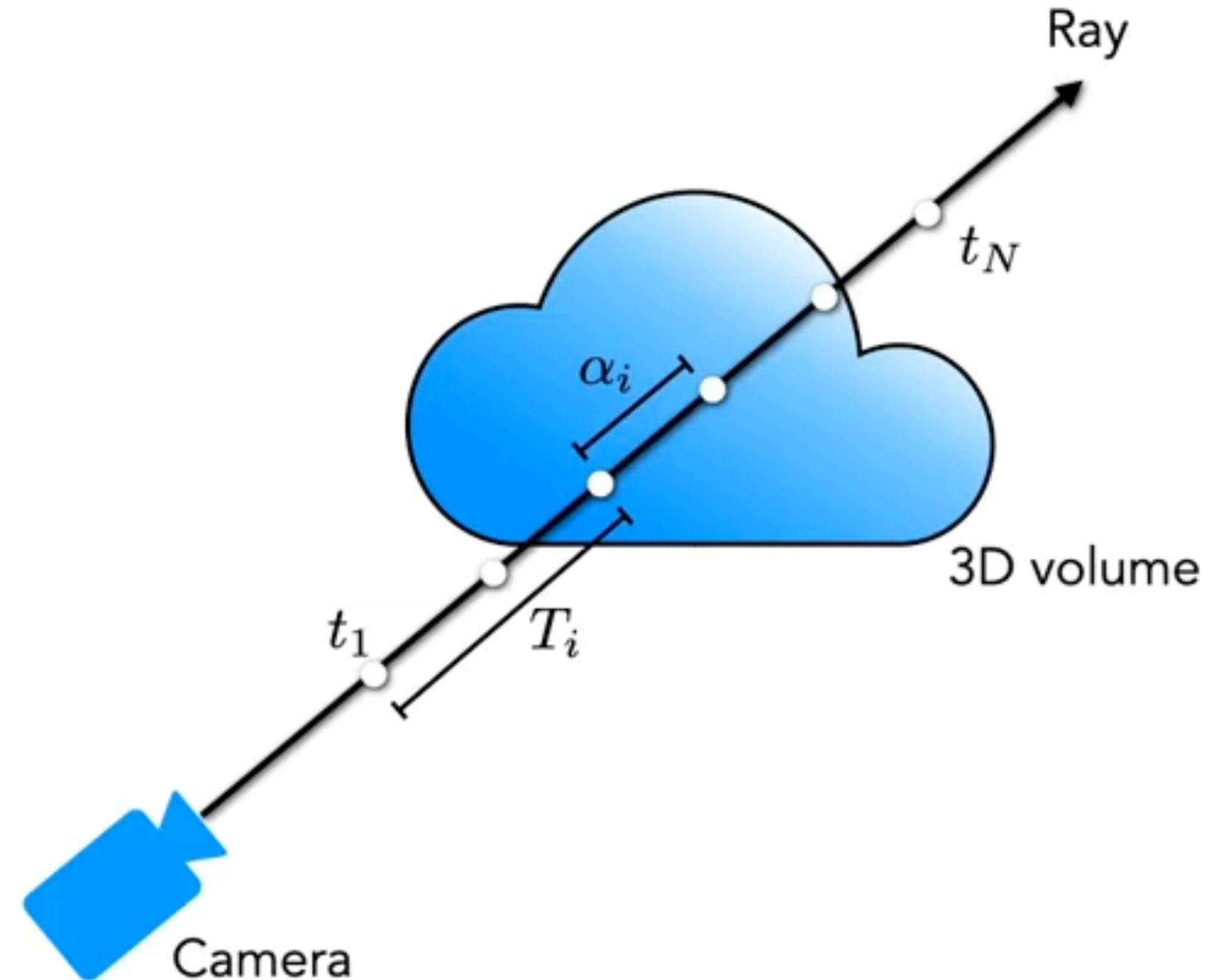
$$C \approx \sum_{i=1}^{N} T_i \alpha_i c_i$$

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

# Numerical integration step easily differentiable

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

colors

weights

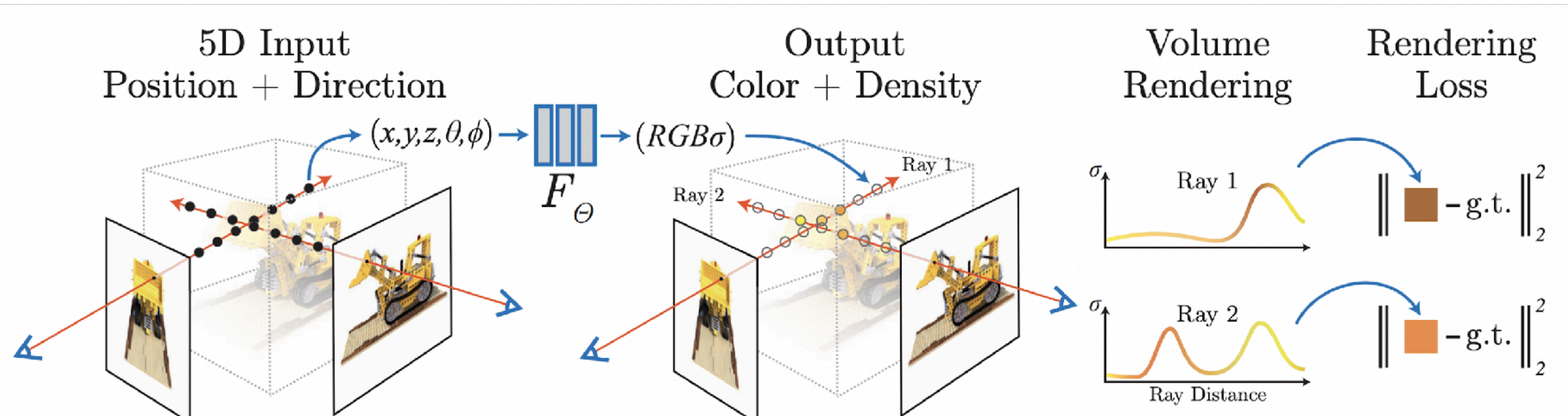**Differentiable w.r.t color $(r, g, b)$ and volumetric density $\sigma$ outputs**

# NeRF Summary

# Results

# Novel View Synthesis



Inputs: sparsely sampled images of scene

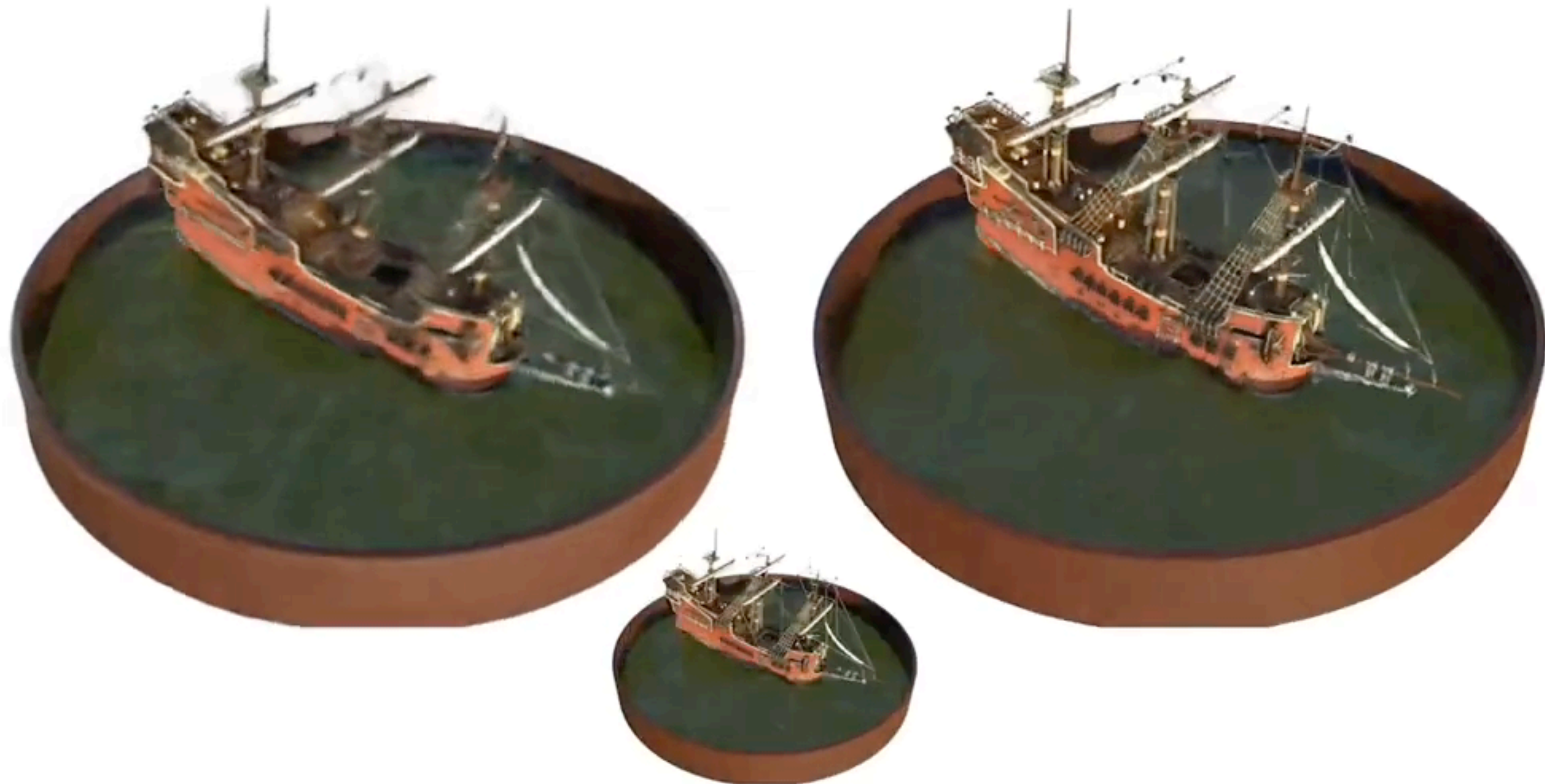Outputs: *new* views of same scene
(rendered by our method)

# More detailed and consistent than prior work that represents scene as discrete voxel grid

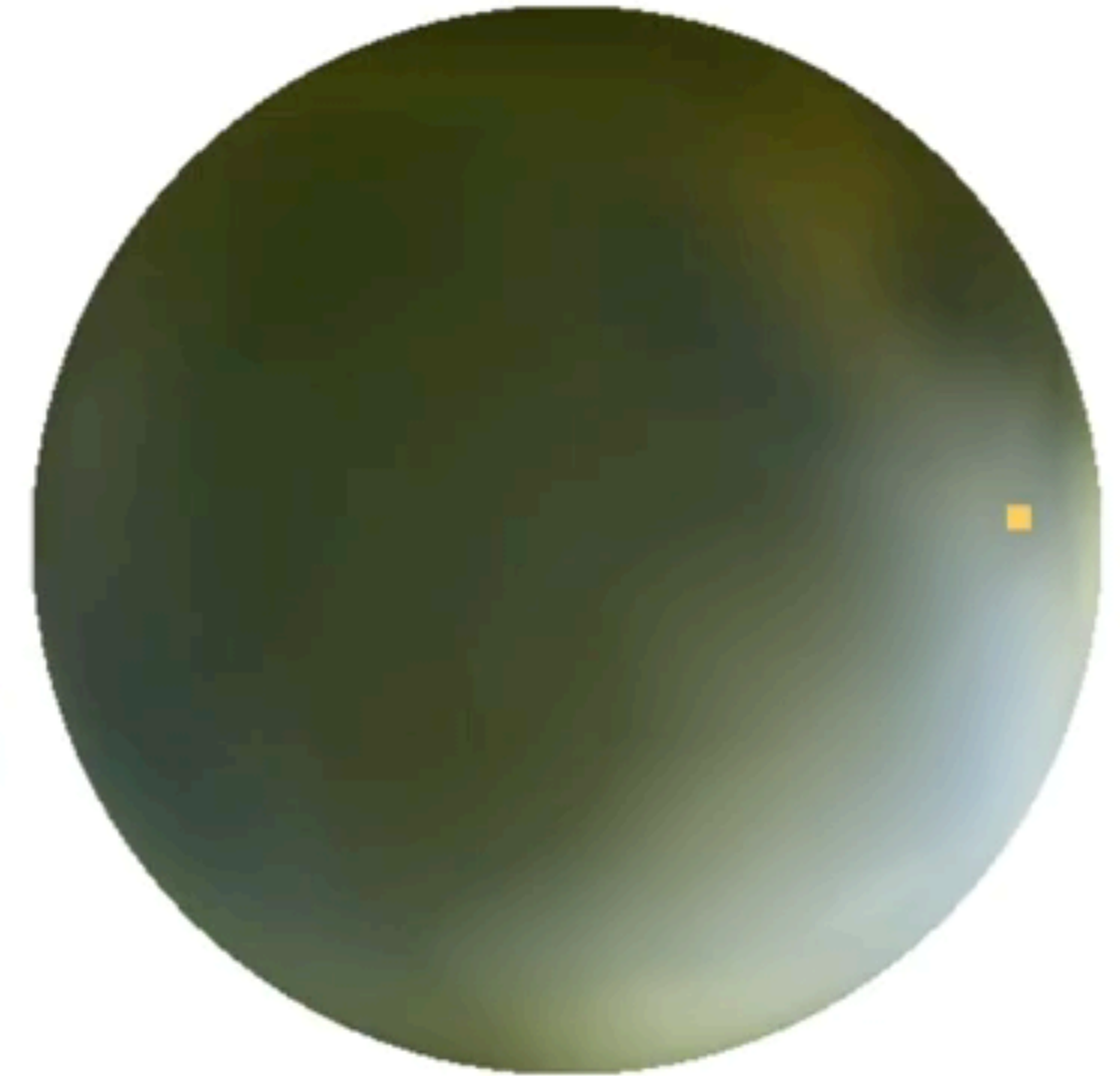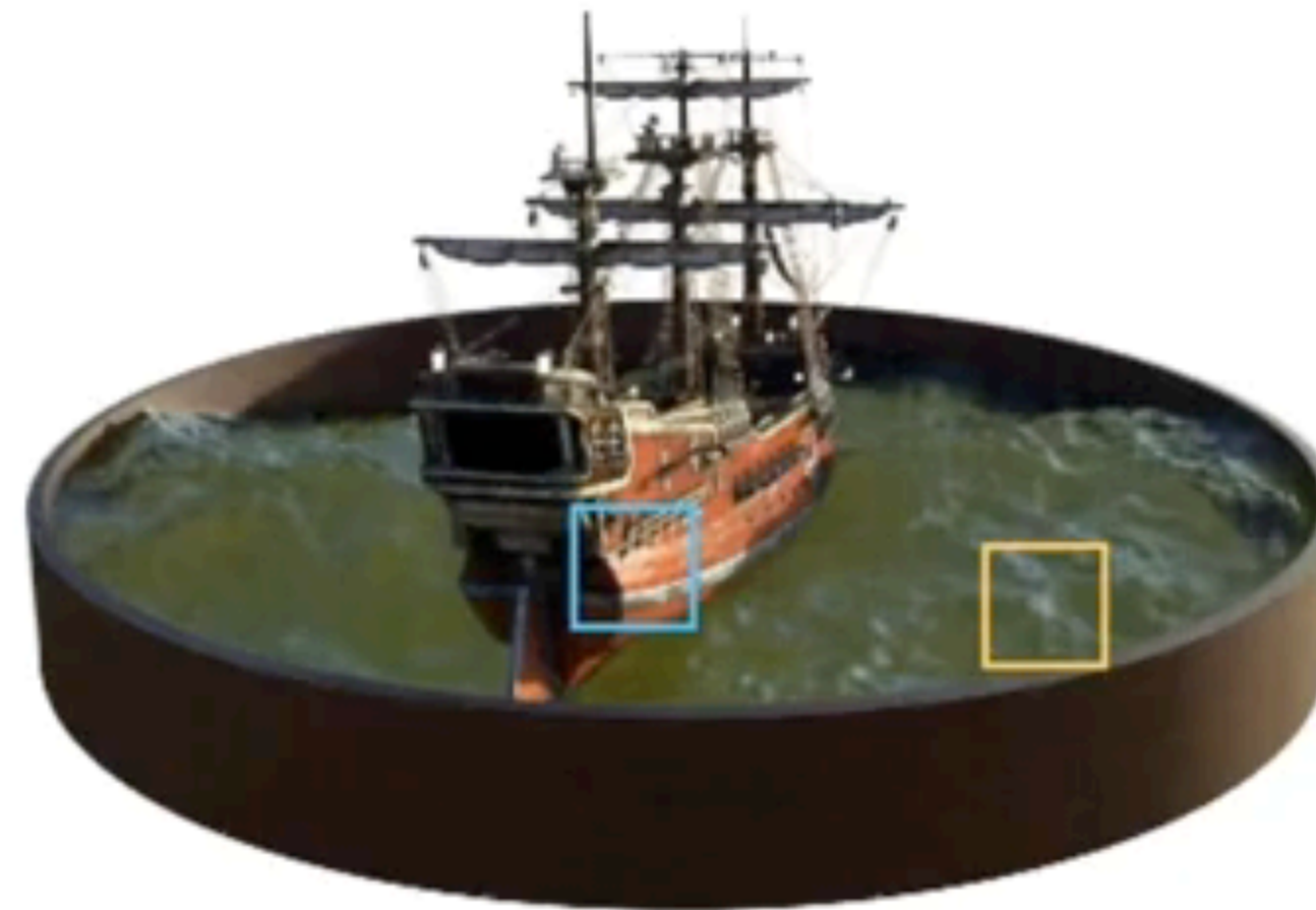Neural Volumes [Lombardi 2019]                                    NeRF



36

Why do we need r,g,b to be a function of viewpoint?

# Viewing directions as input



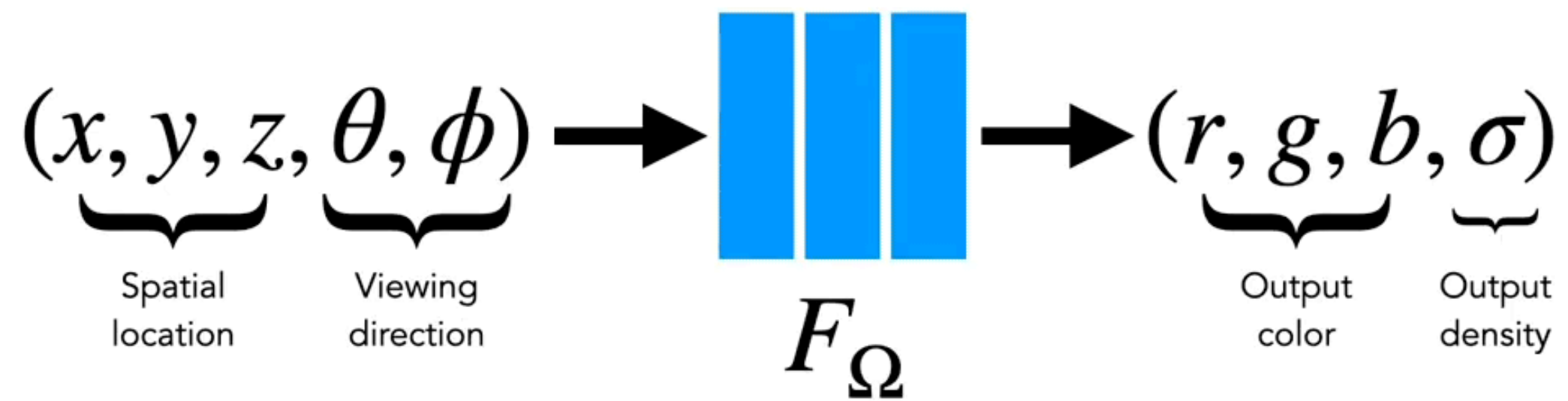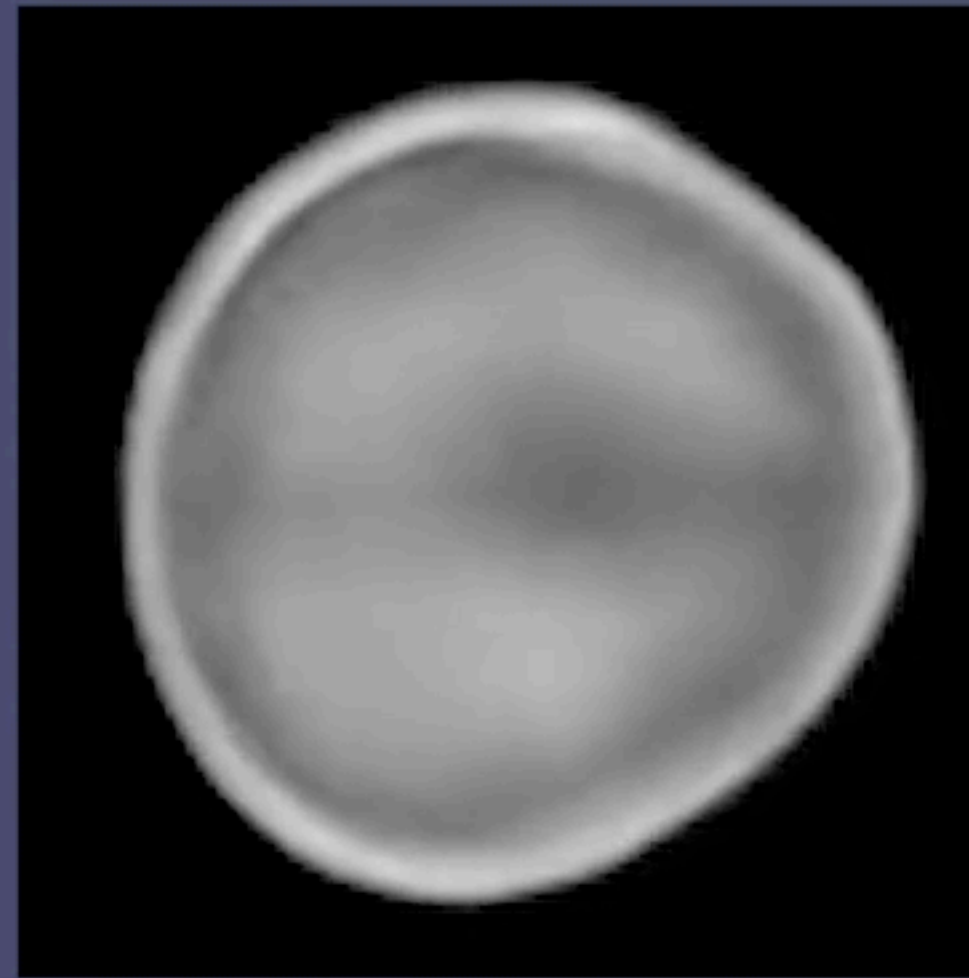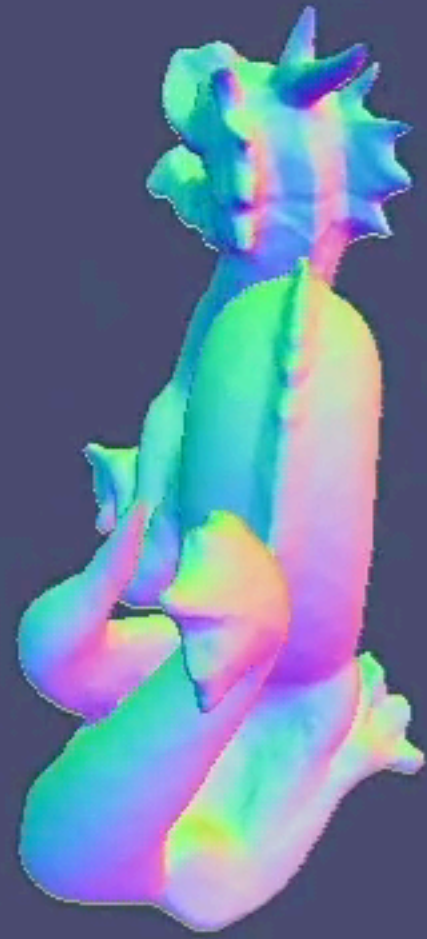Radiance distribution for point on side of ship

Radiance distribution for point on water's surface

One key trick to make it work ...

# Naively passing in position creates blurry images!



Standard input

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

Spatial location  Viewing direction  $F_\Omega$  Output color  Output density

Why?

Let's say we train a
network to memorize an
image

$$(\mathbf{x}) \longrightarrow \blacksquare\blacksquare\blacksquare \longrightarrow (r, g, b)$$

Ground truth image

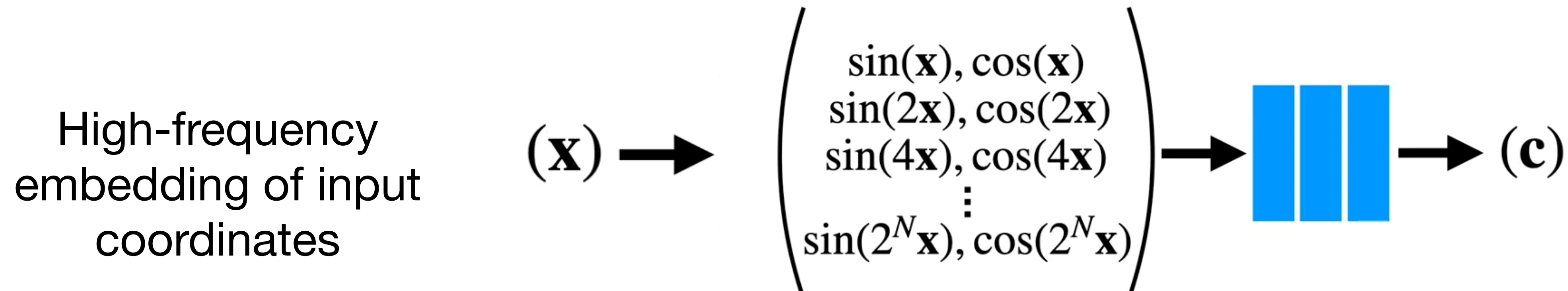Standard fully-connected net

How do we
make the image
look sharper?

# Idea: Encode low-dim coordinates to high-dim features

High-frequency embedding of input coordinates

$$(\mathbf{x}) \longrightarrow \begin{pmatrix} \sin(\mathbf{x}), \cos(\mathbf{x}) \\ \sin(2\mathbf{x}), \cos(2\mathbf{x}) \\ \sin(4\mathbf{x}), \cos(4\mathbf{x}) \\ \vdots \\ \sin(2^N\mathbf{x}), \cos(2^N\mathbf{x}) \end{pmatrix} \longrightarrow \blacksquare \longrightarrow (\mathbf{c})$$
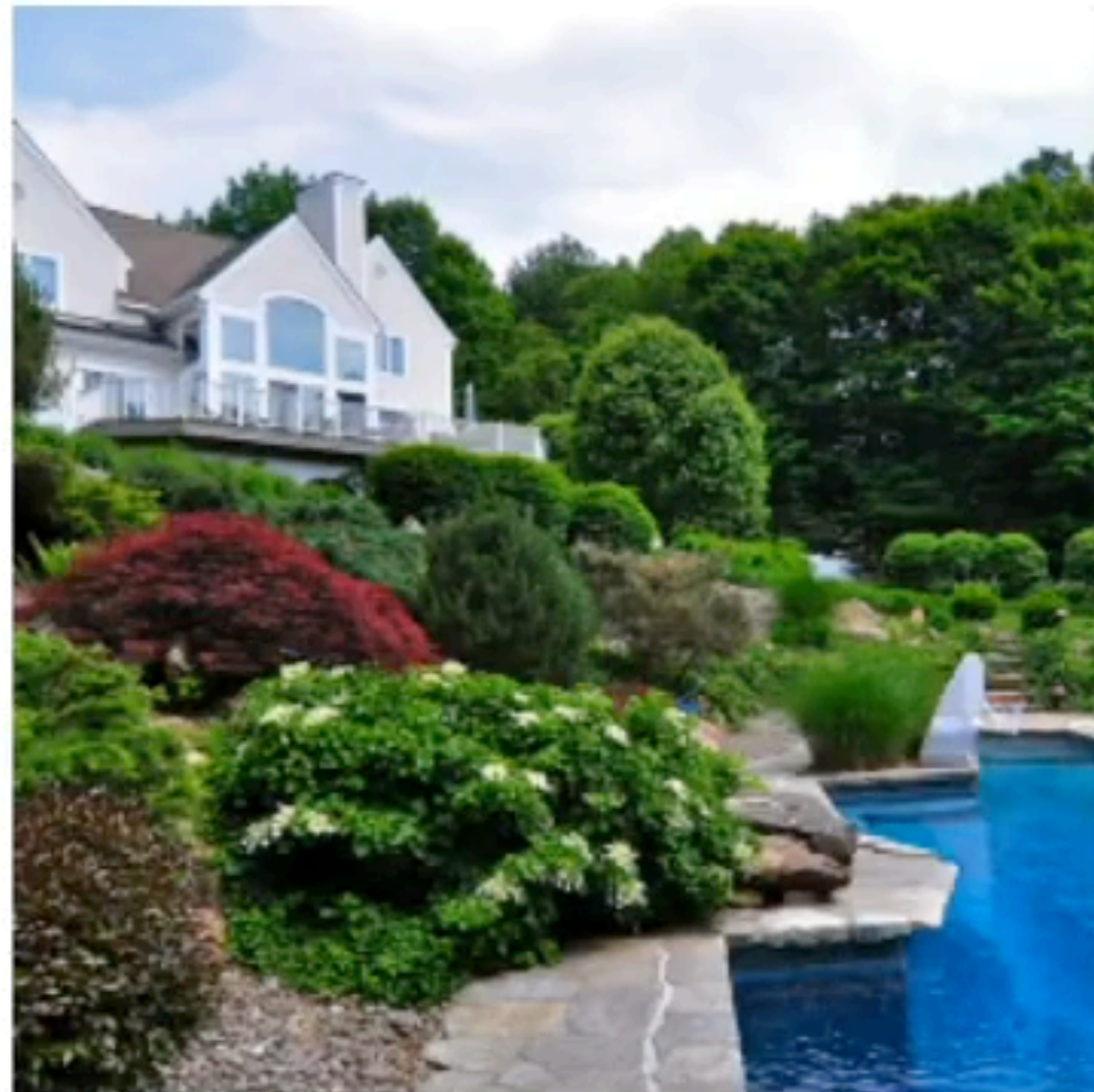
<u>Fourier Features Let Networks Learn High Frequency Functions in Low Dimensional Domains, Tancik et al.</u>

# Use positional encoding

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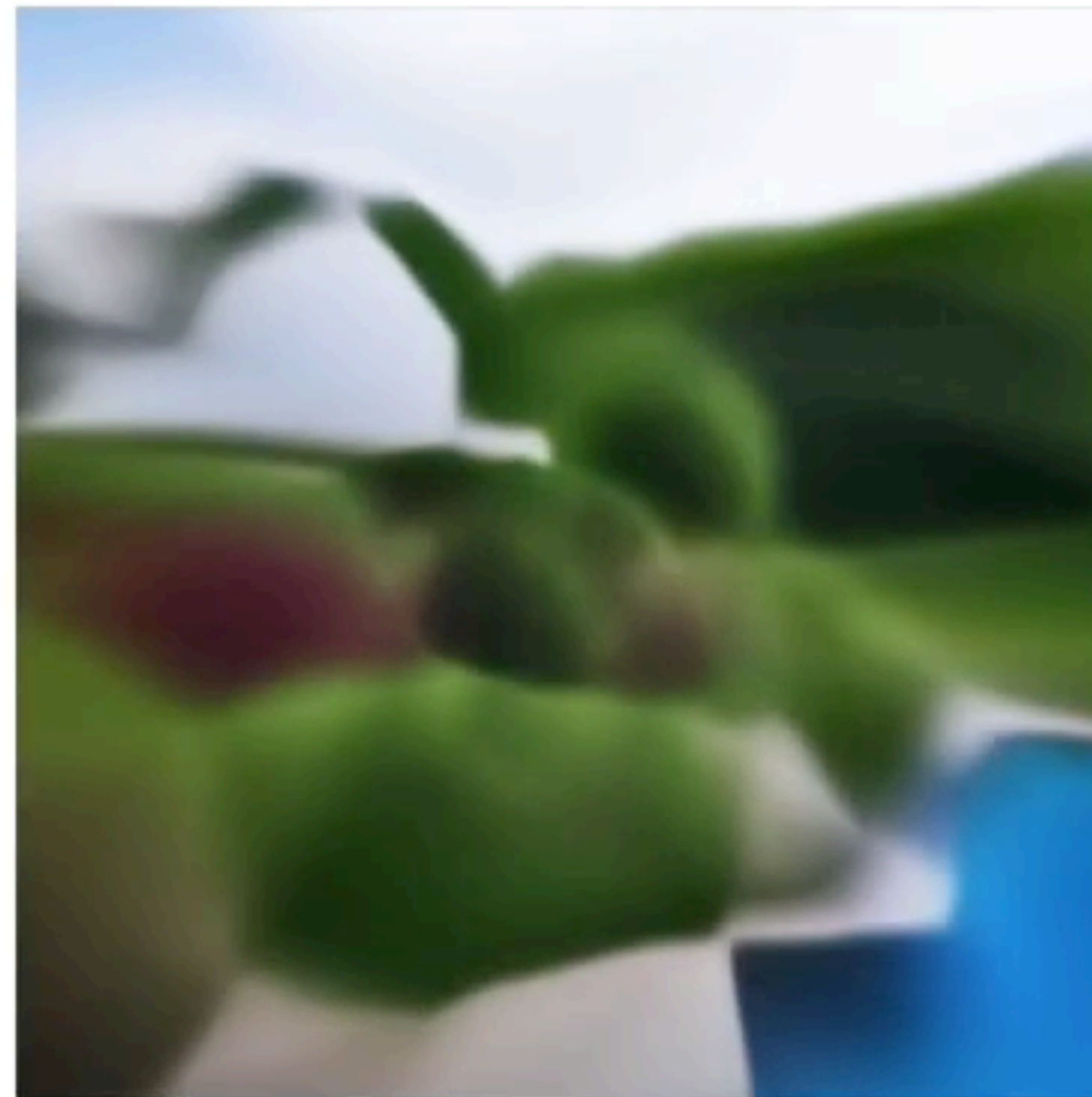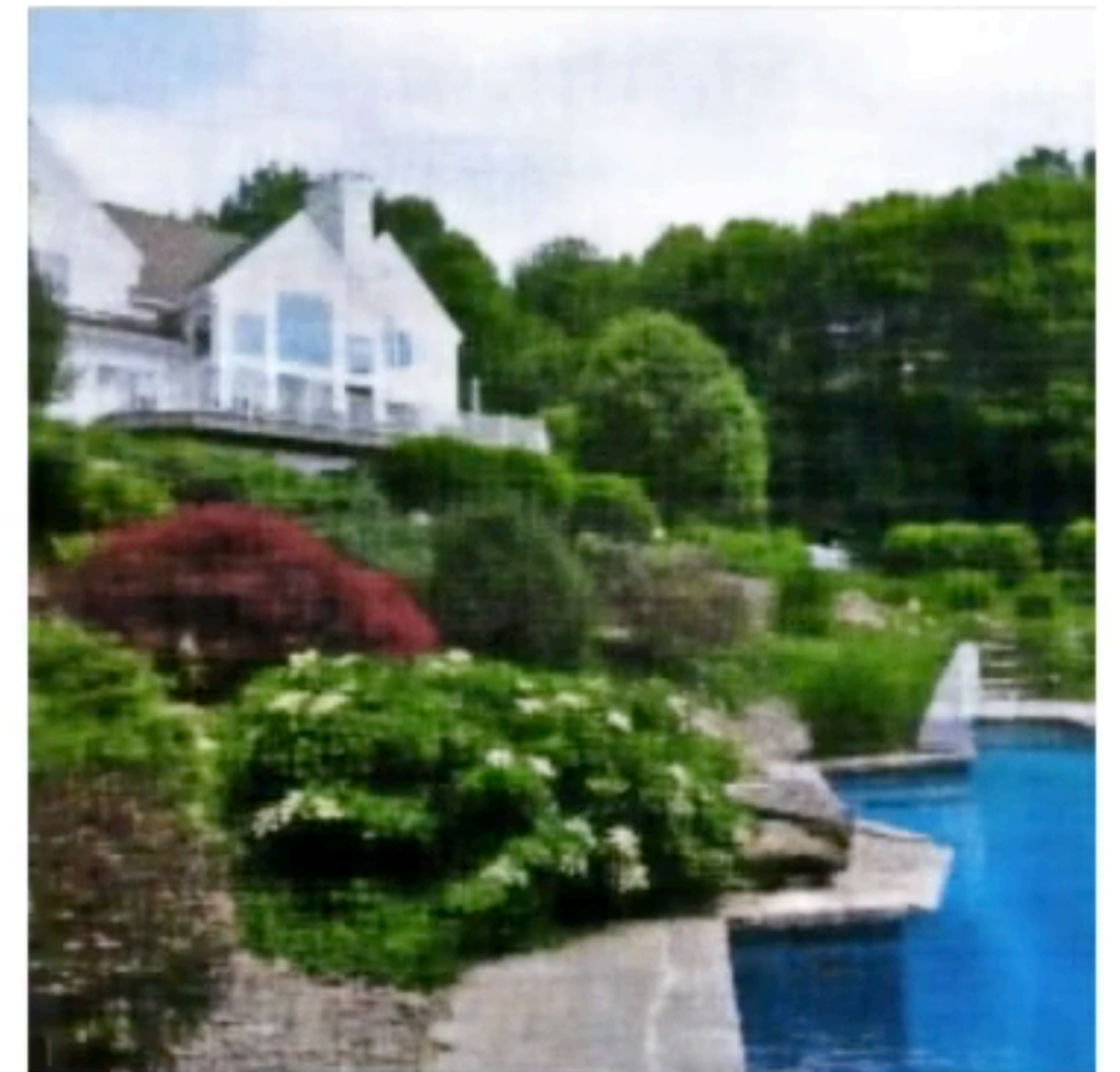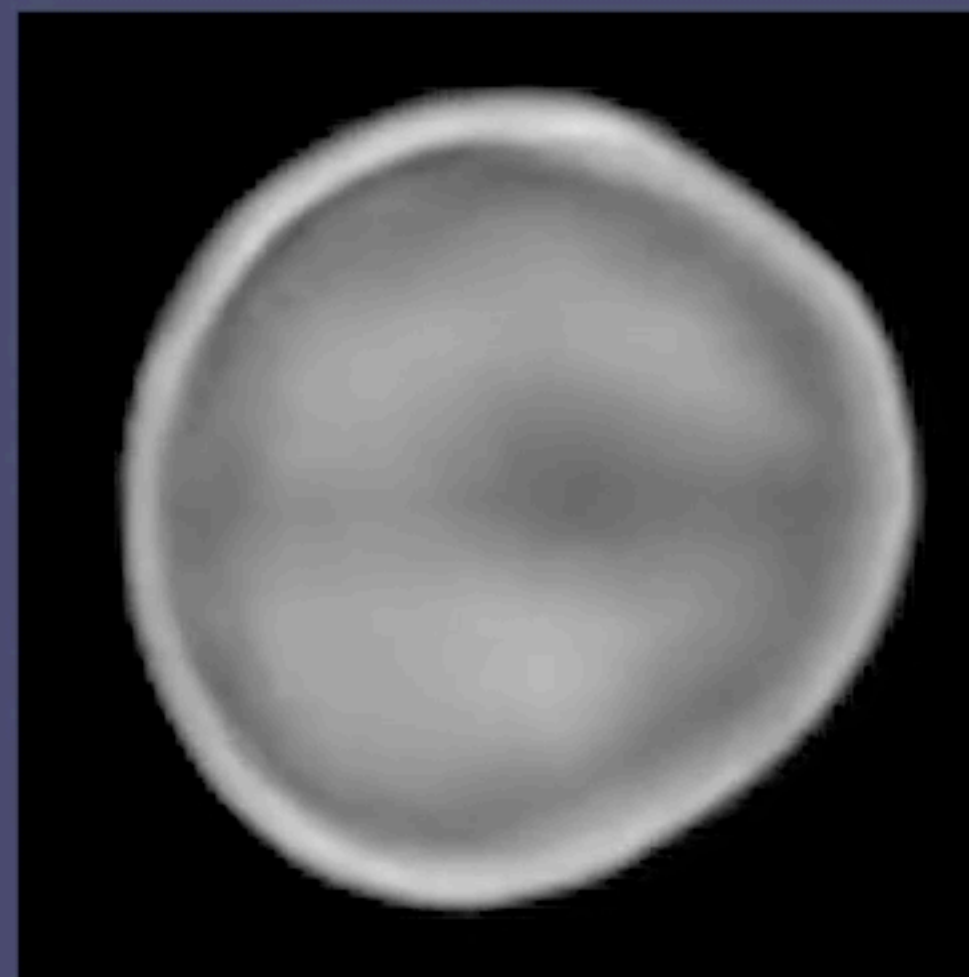
Ground truth image

Standard fully-connected net

With "positional encoding"

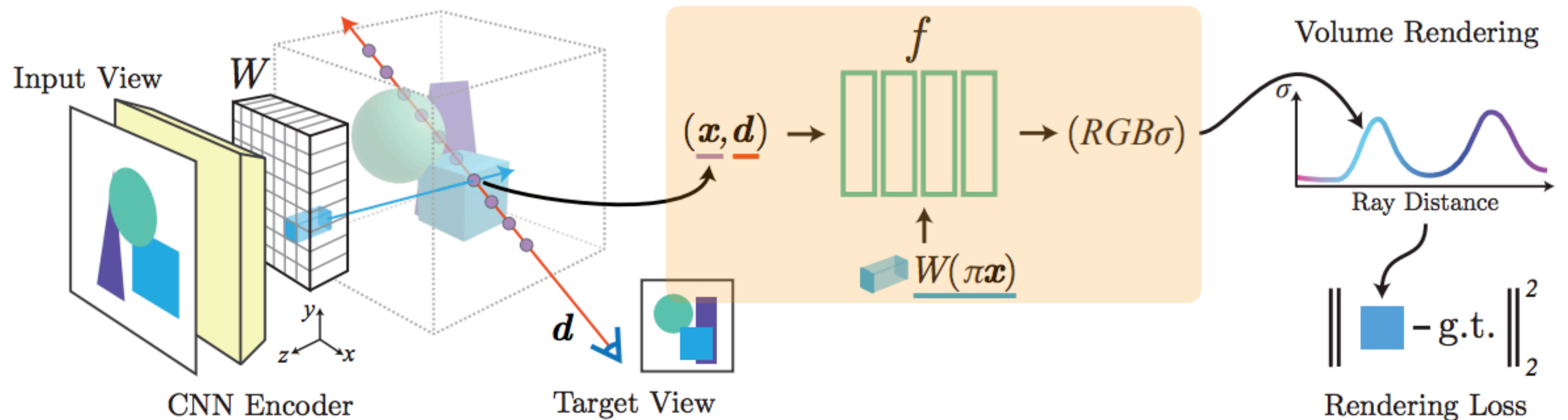Standard input

Fourier feature input

3D Shape

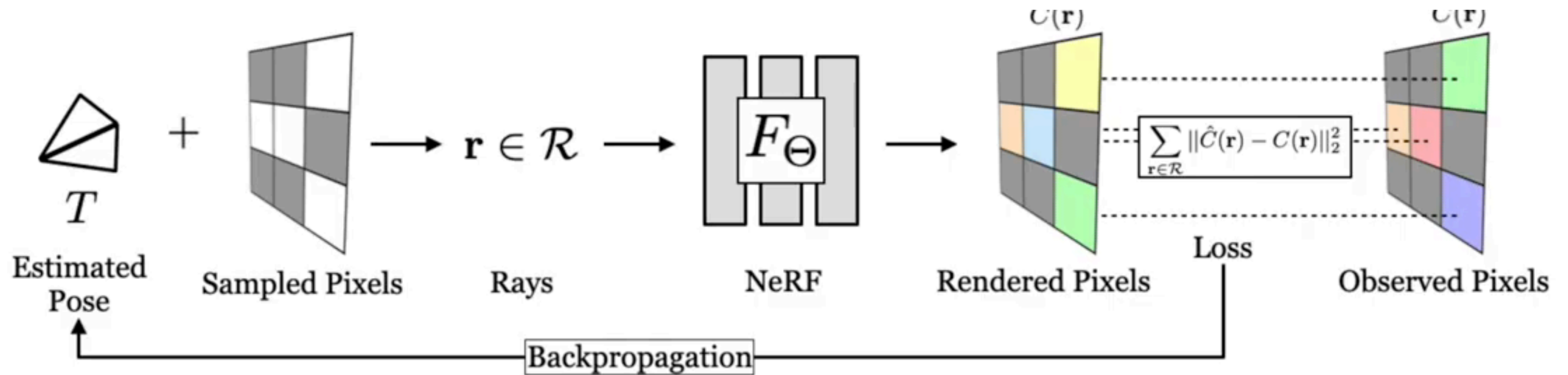3D MRI

3D NeRF

74

# Lots of extension and applications!

# Generalization

Goal: Train a NeRF for arbitrary new scenes with fewer images

# Unknown camera poses

Goal: Estimate poses given a trained NeRF



$$\hat{C}(\mathbf{r})$$

$$\sum_{\mathbf{r}\in\mathcal{R}} \|\hat{C}(\mathbf{r}) - C(\mathbf{r})\|_2^2$$

$$T$$

Estimated Pose     Sampled Pixels    Rays    NeRF    Rendered Pixels    Loss    Observed Pixels

$\mathbf{r} \in \mathcal{R}$

$F_\Theta$

Backpropagation

iNeRF: Inverting Neural Radiance Fields for Pose Estimation

# Unknown camera poses + Scene

Goal: Simultaneously estimate pose and scene representation



NeRF — Images + accurate camera poses → 3D scene representation

BARF (ours) — Images + imperfect camera poses → 3D scene representation + registered camera poses



BARF : Bundle-Adjusting Neural Radiance Fields

iMAP: Implicit Mapping and Positioning in Real-Time

# NERF for Grasping
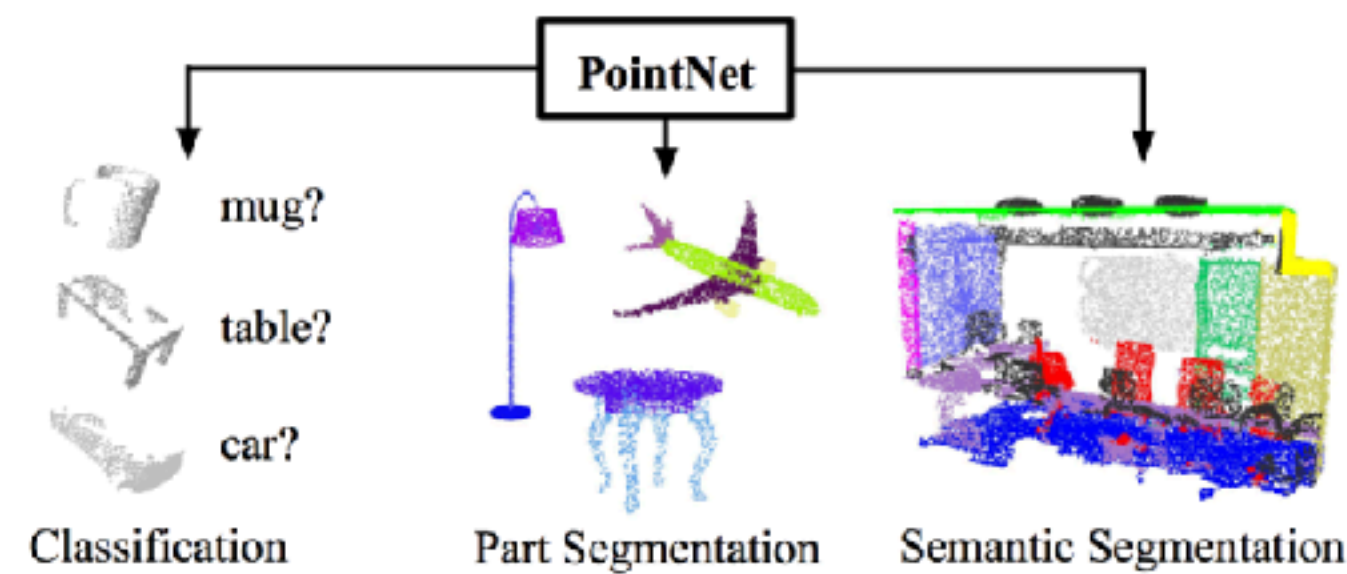


Grasping Singulated Objects

# tl;dr

But manipulating objects require 3D reasoning!

## PointNet

End-to-end learning for **scattered, unordered** point data
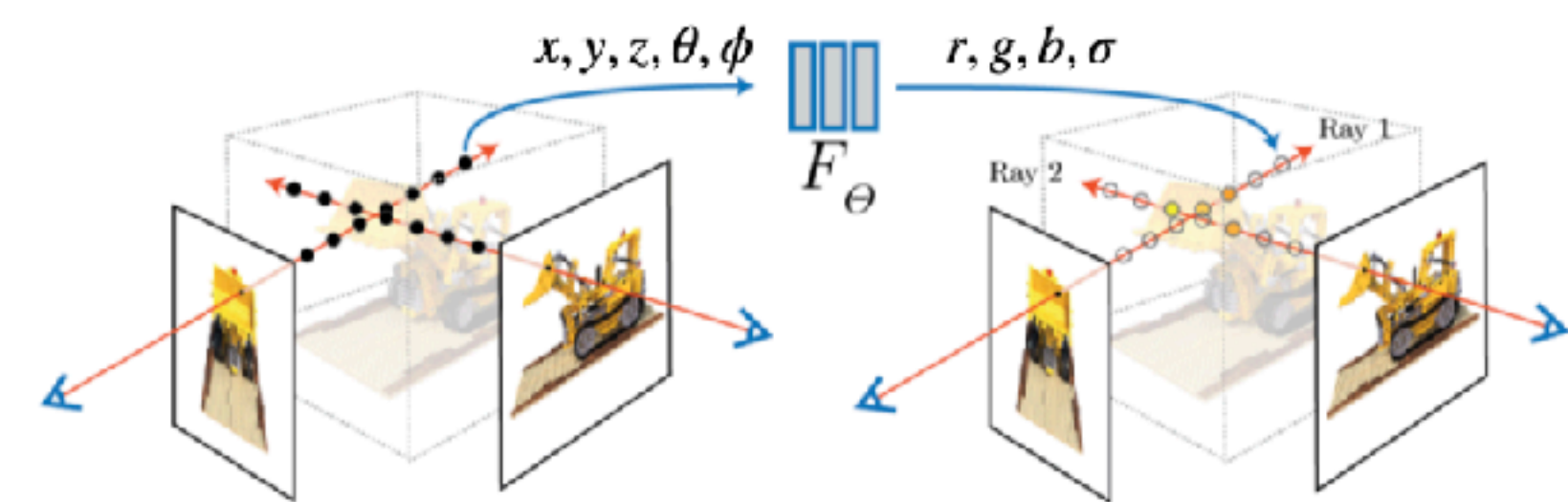
**Unified** framework for various tasks

PointNet

mug?
table?
car?

Classification    Part Segmentation    Semantic Segmentation

Slides from Qi et al, CVP 2017  http://stanford.edu/~rqi/pointnet/docs/cvpr17_pointnet_slides.pdf

13

## What are Neural Radiance Fields (NeRFs)?

Idea: Use a neural network to *implicitly* represent 3D volume!

$x, y, z, \theta, \phi$    $F_\Theta$    $r, g, b, \sigma$

Ray 1
Ray 2

✓ No Discretization    ✓ Compressible    ✓ Differentiable

54