# COS 426: Computer Graphics

# Neural Rendering

Felix Heide

PRINCETON UNIVERSITY

# .... so far so good: Computer Graphics

# Rendering Equation

$$L_{\mathrm{o}}(\mathbf{p}, \omega_{\mathrm{o}}) = L_{\mathrm{e}}(\mathbf{p}, \omega_{\mathrm{o}}) + \int_{\Omega} L_{\mathrm{i}}(\mathbf{p}, \omega_{\mathrm{i}}) f_{\mathrm{r}}(\mathbf{p}, \omega_{\mathrm{i}}, \omega_{\mathrm{o}}) (\omega_{\mathrm{i}} \cdot \mathbf{n}) \, \mathrm{d}\omega_{\mathrm{i}}$$

Outgoing radiance        Incident radiance    BRDF

# Rendering Equation

$$L_{\mathrm{o}}(\mathbf{p}, \omega_{\mathrm{o}}) = L_{\mathrm{e}}(\mathbf{p}, \omega_{\mathrm{o}}) + \int_{\Omega} \boxed{L_{\mathrm{i}}(\mathbf{p}, \omega_{\mathrm{i}})} f_{\mathrm{r}}(\mathbf{p}, \omega_{\mathrm{i}}, \omega_{\mathrm{o}}) \left(\omega_{\mathrm{i}} \cdot \mathbf{n}\right) \mathrm{d}\,\omega_{\mathrm{i}}$$

# Generative Adversarial Networks



Prior Distribution → Sample → Generator → Synthesis → Discriminator → Classification Real/Synthesized
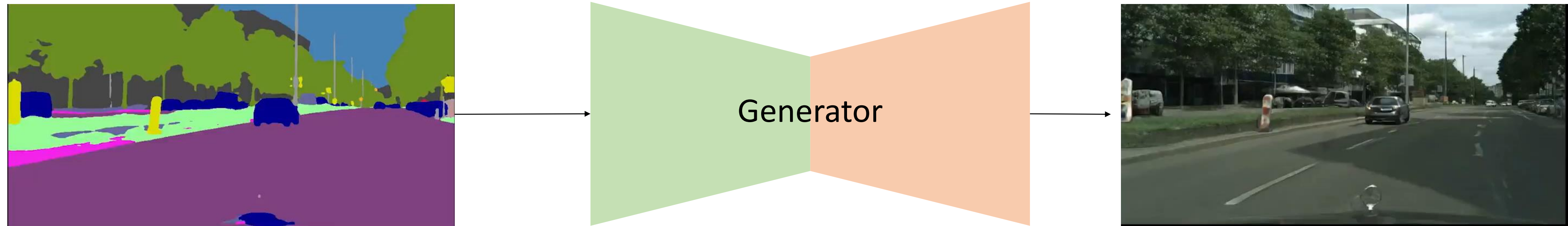
Real images

# Generative Adversarial Networks



StyleGAN [Karras et al., 2019]

# Conditional Generative Models



Generator

Vid2Vid [Wang et al., 2019]
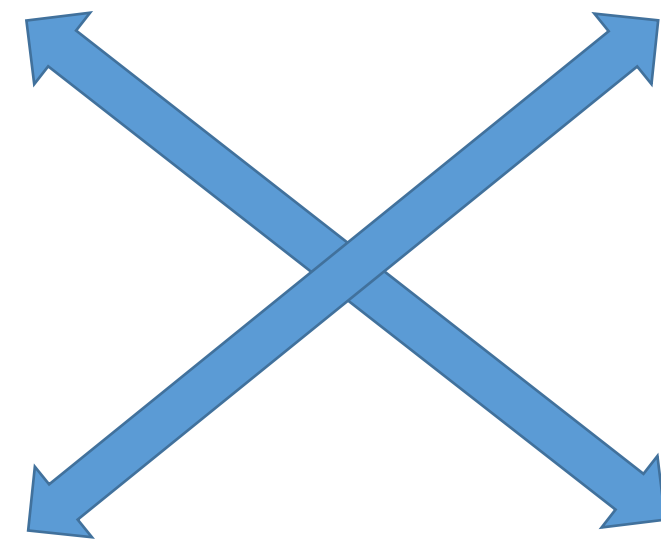
# Two Alternatives of Realistic Image Synthesis

## Photo-realistic Rendering (CG)

### Cons:

- Requires lots of manual work
  - Building of high-quality assets
  - Setting up the scene
- Long render times

### Pros:

- Full control of scene parameter:
  - Camera, light sources, motion, geometry, appearance

## Generative Machine Learning (ML)

### Cons:

- Requires lots of training data
- No fine-grained semantic control of the scene parameters, e.g., motion or illumination
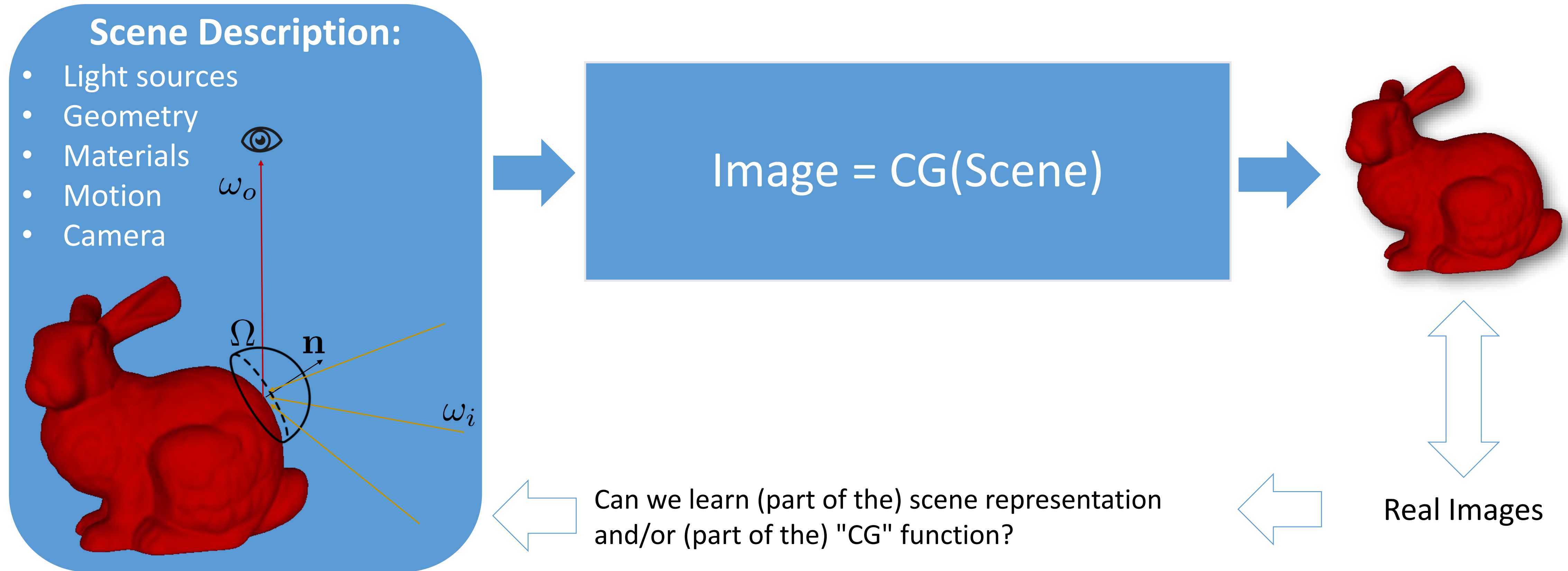
### Pros:

- Fully automatic training
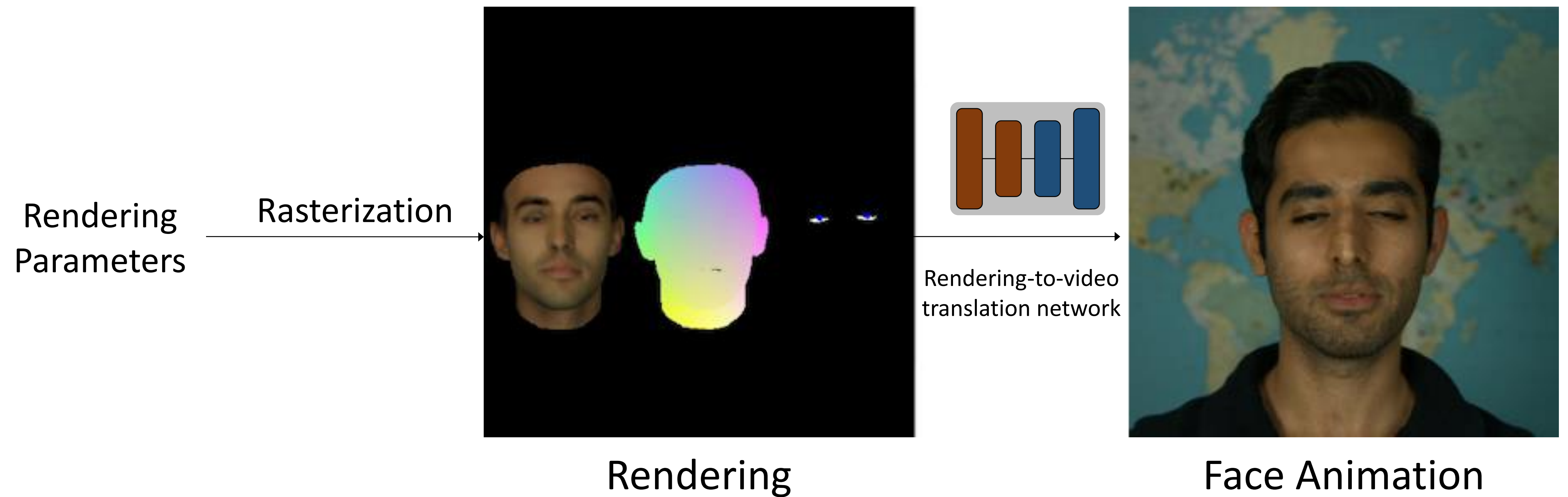- Interactive inference/rendering

**Neural** Rendering to the rescue!

# Neural Rendering - Graphics vs. Learning



**Scene Description:**
- Light sources
- Geometry
- Materials
- Motion
- Camera

$\omega_o$

$\Omega$

$\mathbf{n}$

$\omega_i$

Image = CG(Scene)

Real Images

Can we learn (part of the) scene representation and/or (part of the) "CG" function?

Neural Rendering to the rescue!

# CG Modules

Computer Graphics Module: Rasterization is used to synthesize the input to the network.



Rendering
Parameters → Rasterization → **Rendering** → Rendering-to-video translation network → **Face Animation**

Deep Video Portraits [Kim et al., 2018]

# CG Modules

Computer Graphics Module: *Differentiable* Volume Renderer



Neural Volumes [Lombardi et al., 2019]

# Neural Rendering Definition:

> "(Deep) neural networks for image or video generation that enable explicit or implicit control of scene properties"



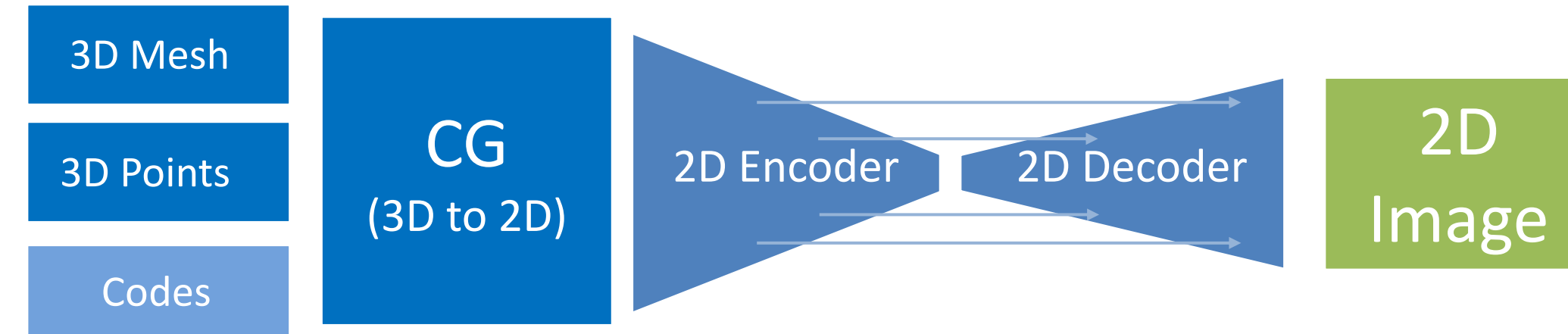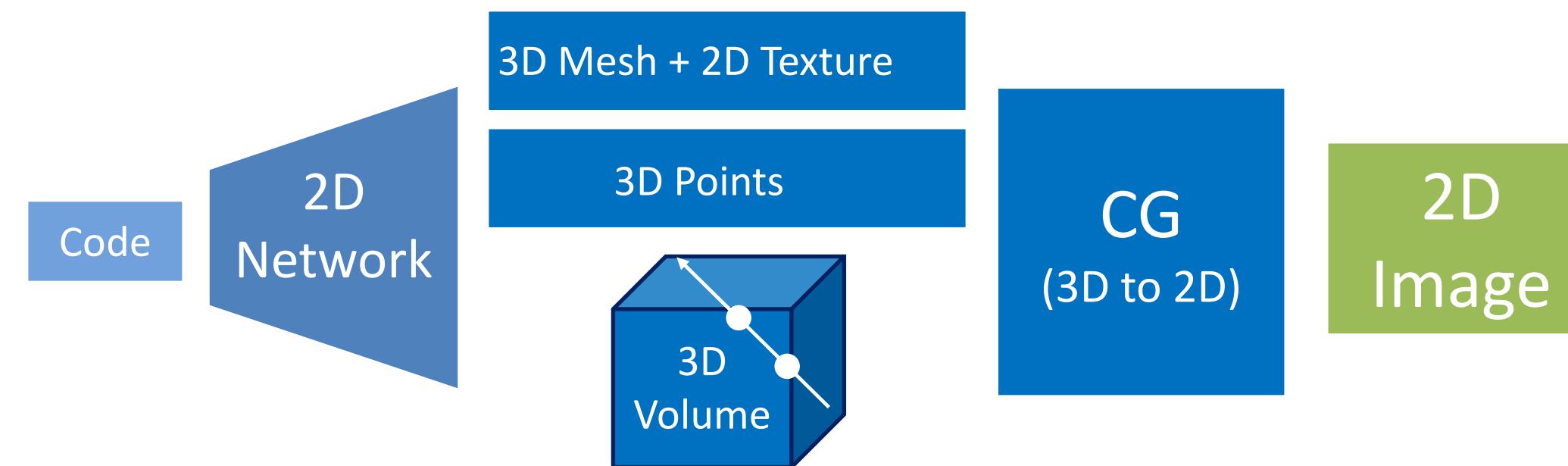Allows us to create photorealistic assets
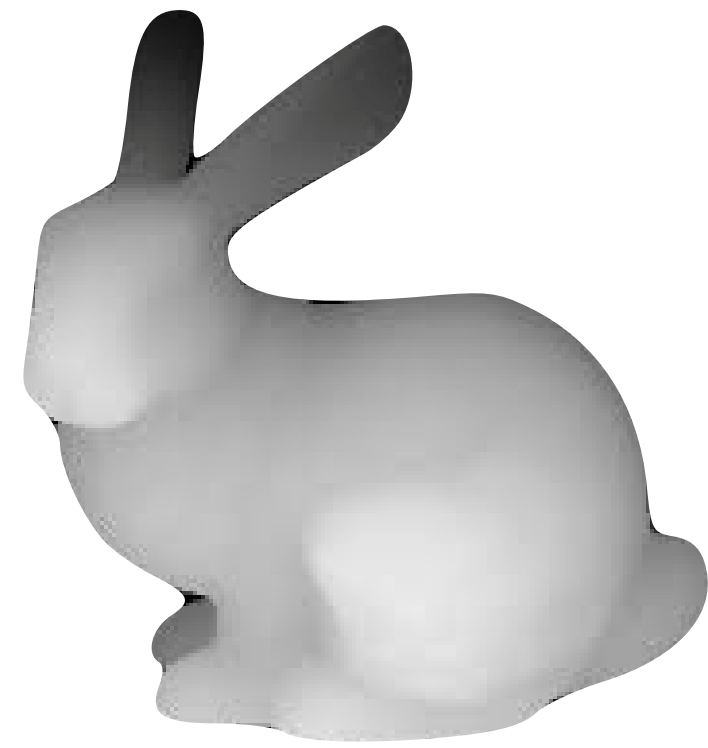
# Neural Rendering Zoo

# What's the right scene representation?



**Point cloud**

**Depth map**

**Voxel grid**

**Triangle mesh**

# What's the right scene representation?
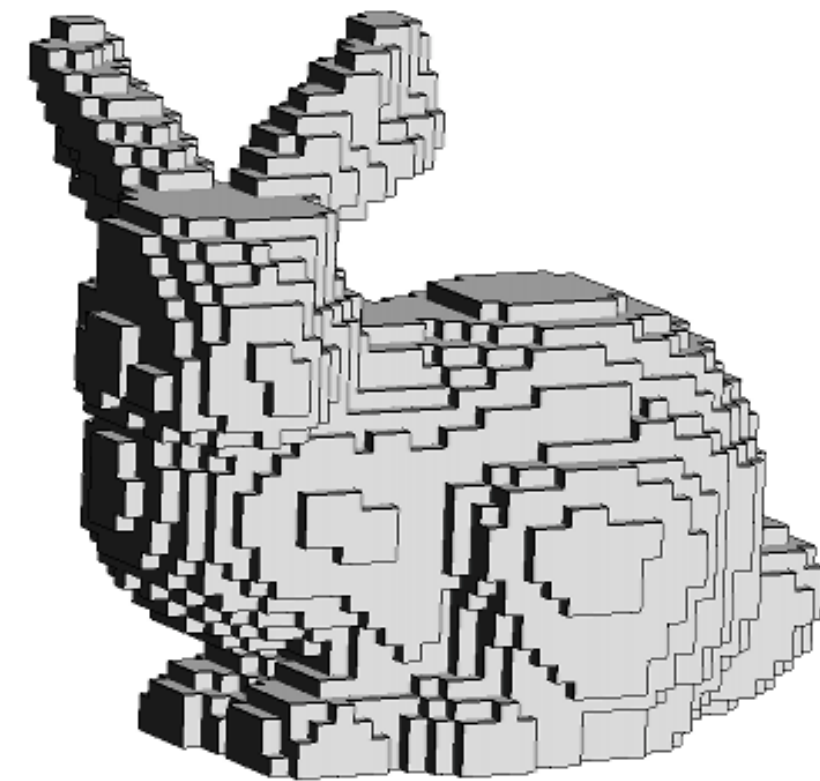
▸ **Do sensors output data in this representation?**



**Point cloud**          Depth map          Voxel grid          Triangle mesh

# What's the right scene representation?

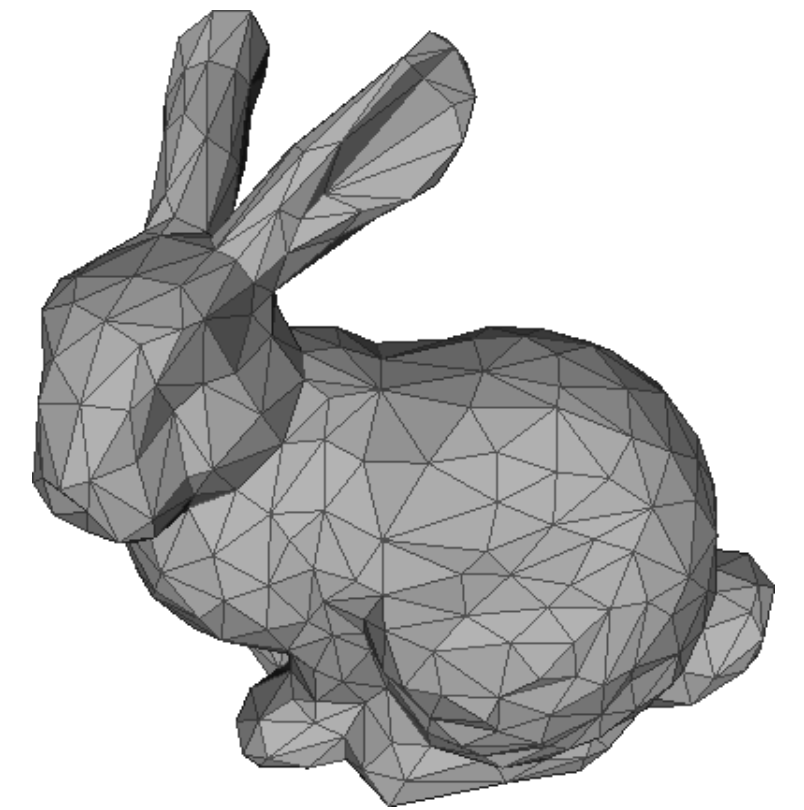▸ **Can we process/generate content in this representation?**



Point cloud

**Depth map**

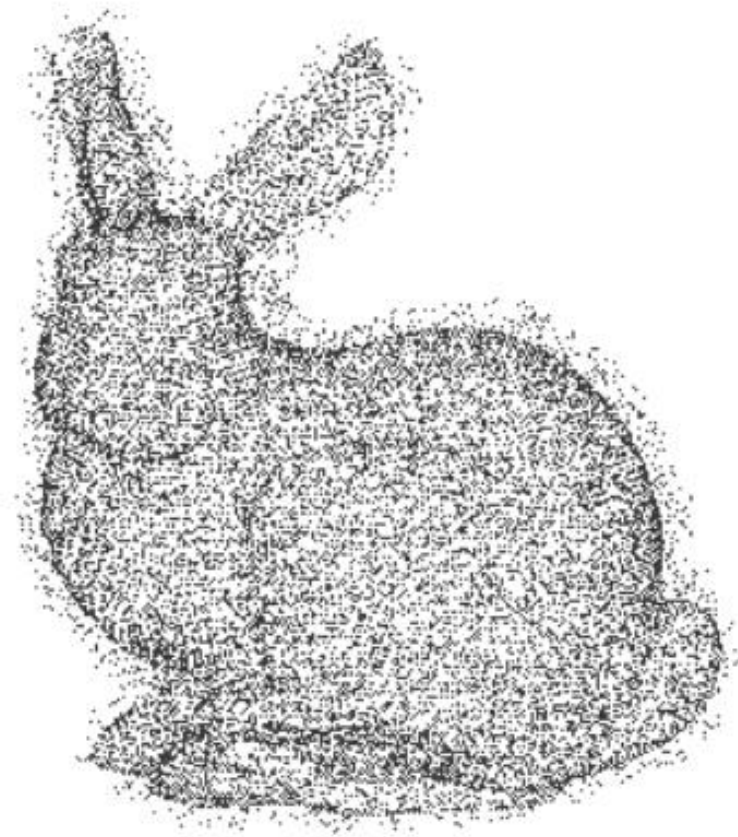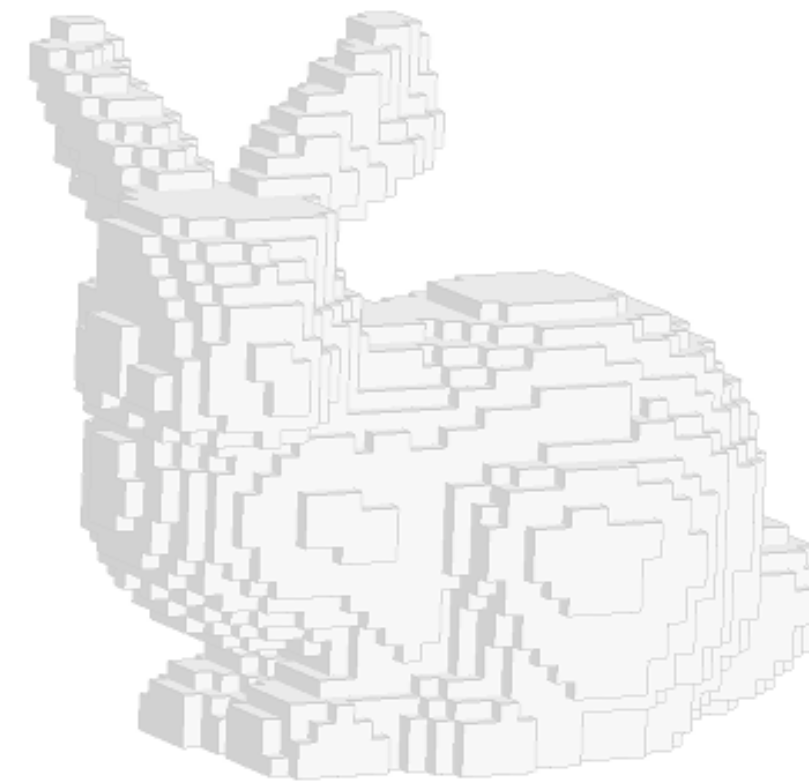**Voxel grid**

Triangle mesh

# What's the right scene representation?

▸ Do sensors output data in this representation?

▸ Can we process/generate content in this representation?

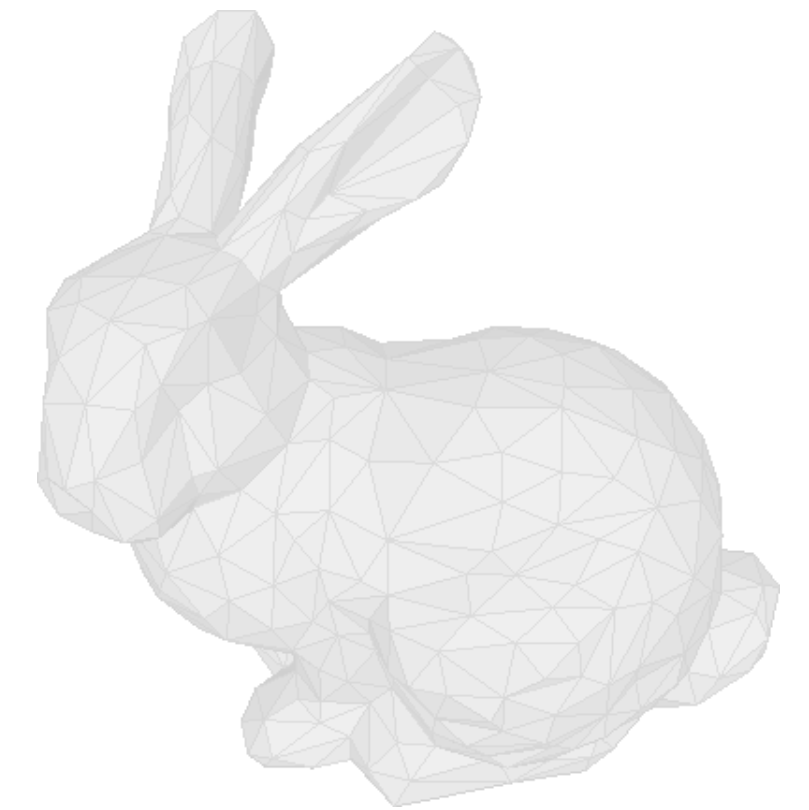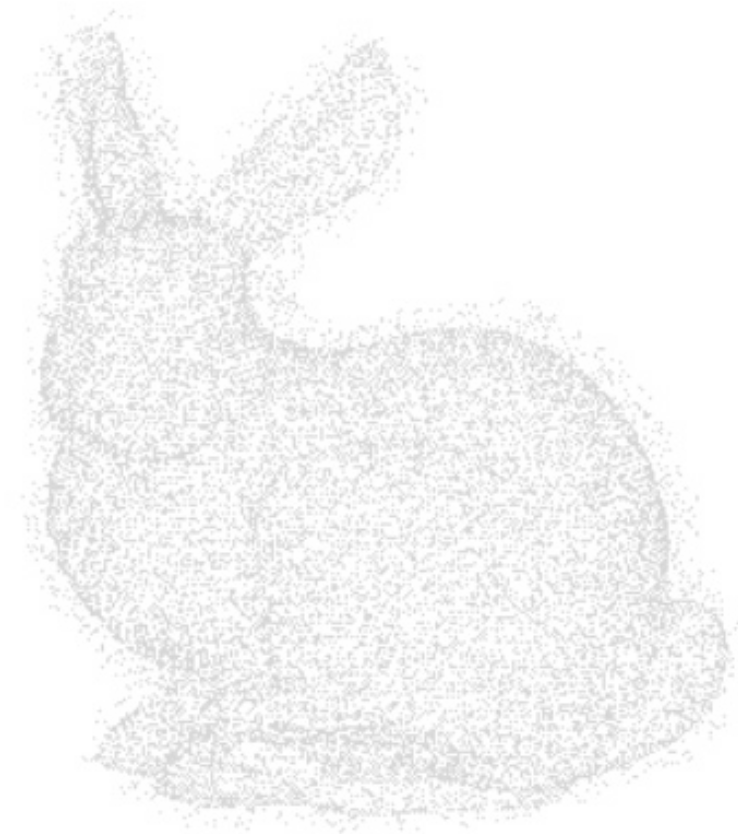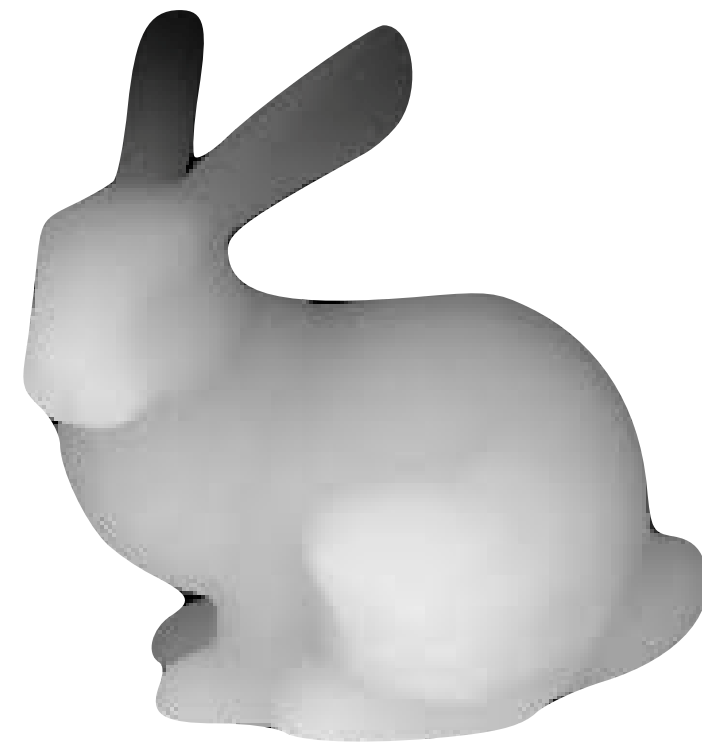▸ **Can we easily render this representation?**



Point cloud

Depth map

Voxel grid

**Triangle mesh**

# What's the right scene representation for 3D vision?

▸ Sensors don't give us nice data

▸ World is not dense in 3D space, so maybe representation shouldn't be

▸ Doing computation directly on "efficient" representations is hard

▸ What's missing? Efficient representation that's easy to optimize with gradient-based methods!

# NeRF (neural radiance fields):

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

# Neural Volumetric Rendering

# Neural Volumetric Rendering

What's the radiance/color arriving at this pixel?

▶ "Soft" volumetric functions better suited for gradient-based optimization

# Neural Volumetric Rendering

▶ **(Coordinate-based) neural network represents scene as continuous function**

# NeRF: neural volumetric rendering for view synthesis



Inputs: sparsely sampled images of scene

Outputs: new views of same scene

NeRF in the Wild, Martin-Brualla et al.

NeRF in the Wild, M



NeRFies, Park et al.

NeRF in the Wild

NeRF

Neural Scene Flow Fields, Li et al.

NeRF in the Wild

NeRF

Neural Scene Flow Fields, Li et al.

Dynamic Neural Radiance Fields, Gafni et al.

Control Pose

Control Expression

NeRF in the Wild

Neural Scene Graphs, Ost et al.

NeRF

Neural Scene Flow Fields, Li et al.

Control Pose

Control Expression

Dynamic Neural Radiance Fields, Gafni et al.

# Representing a scene as a continuous 5D function

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

Spatial location

Viewing direction

$F_\Omega$

Fully-connected neural network
9 layers,
256 channels

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

Output color

Output density

# Neural network replaces large N-d array: tradeoff between storage and computation



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

versus

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

# Generate views with traditional volume rendering

# Generate views with traditional volume rendering

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

# 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

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

$t_N$

$t_1$

$T_i$

3D volume

Camera

# Generate views with traditional volume rendering

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

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

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*:

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



Ray

$t_N$

$\alpha_i$

3D volume

$t_1$    $T_i$

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



Radiance distribution for point on side of ship

Radiance distribution for point on water's surface

# 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



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

# Optimize with gradient descent on rendering loss

Any differentiable scene representation $F_\Omega$
could be used here

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

# 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)

# Challenge:

How to get MLPs to represent higher frequency functions?

# Simpler toy problem: memorizing a 2D image

$$(x, y) \rightarrow \blacksquare\blacksquare\blacksquare \rightarrow (r, g, b)$$

# Simple trick enables network to memorize images

Ground truth image

Standard fully-connected net

# Positional encoding

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

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

# Training networks ≈ kernel regression

▶ Recent ML theory work shows that training neural network with gradient descent becomes the same as performing kernel regression as the width of each layer goes to infinity

▶ Can examine corresponding kernel function (the neural tangent kernel) to see why adding Fourier feature mapping allows MLPs to represent high frequency functions

Jacot et al., Neural Tangent Kernel: Convergence and generalization in neural networks, NeurIPS 2018

# Kernel regression

▶ Method for fitting a continuous function to a set of data points $\{(x_i, y_i)\}$

▶ High level: add up a set of blobs (kernel functions), one centered at each input point, each with its own weight

▶ Weights are optimal in a least-squares sense: $\min_w \sum_i \| y_i - \hat{f}_w(x_i) \|^2$

$$\hat{f}_w(x) = \sum_{i=1}^{n} w_i k(x - x_i)$$

Blob centered at training input point $x_i$

Estimated function

Weight corresponding to blob centered at $x_i$

# "Width" of kernel function is critical

▶ If the kernel function is too wide, reconstruction is too smooth. If it's too skinny, reconstruction does not interpolate correctly.

▶ Similar to picking the right reconstruction filter bandwidth in signal processing to avoid either blurring or aliasing.

# "Width" of kernel function is critical



Nonparametric function estimate

Kernel shape

Avg distance between training points

Underlying kernel function

Training data
(all yellow points)

Function estimated by
kernel regression

# "Width" of kernel function is critical

# Training networks ≈ kernel regression

▶ Recent ML theory work shows that training neural network with gradient descent becomes the same as performing kernel regression as the width of each layer goes to infinity

▶ Using a Fourier feature mapping changes the corresponding kernel function (the neural tangent kernel), allowing MLPs to represent higher frequency functions

Jacot et al., Neural Tangent Kernel: Convergence and generalization in neural networks, NeurIPS 2018

# Fourier feature mapping: simple 2D example

$$\gamma(\mathbf{X})$$

$\sin(2\pi\mathbf{b}_1^\top\mathbf{x})$  $\cos(2\pi\mathbf{b}_1^\top\mathbf{x})$

$\sin(2\pi\mathbf{b}_2^\top\mathbf{x})$  $\cos(2\pi\mathbf{b}_2^\top\mathbf{x})$

$\sin(2\pi\mathbf{b}_3^\top\mathbf{x})$  $\cos(2\pi\mathbf{b}_3^\top\mathbf{x})$

$\vdots$

$\mathbf{X}$

$x_1$

$x_2$

$\sin(2\pi\mathbf{b}_{N-2}^\top\mathbf{x})$  $\cos(2\pi\mathbf{b}_{N-2}^\top\mathbf{x})$

$\sin(2\pi\mathbf{b}_{N-1}^\top\mathbf{x})$  $\cos(2\pi\mathbf{b}_{N-1}^\top\mathbf{x})$

$\sin(2\pi\mathbf{b}_N^\top\mathbf{x})$  $\cos(2\pi\mathbf{b}_N^\top\mathbf{x})$

# Simple trick enables network to memorize images

Ground truth image

Standard fully-connected net

With "encoding"

# Positional encoding also directly improves our scene representation!
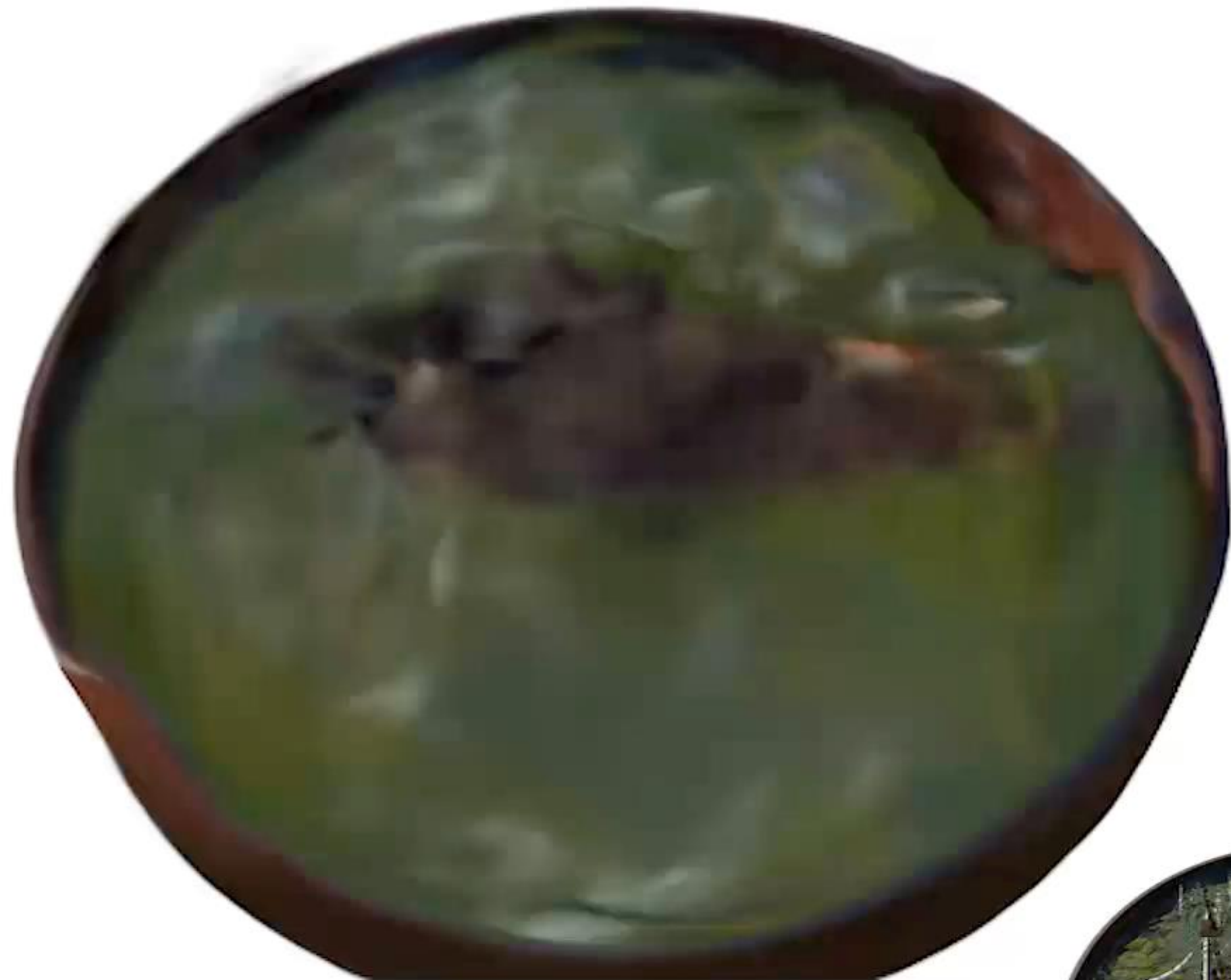


NeRF (Naive)

NeRF (with positional encoding)

Check out https://bmild.github.io/fourfeat/ for more details

# More detailed and consistent than prior work that represents scene as function encoded by MLP

SRN [Sitzmann et al. 2019]

NeRF



Nearest Input

# NeRF encodes convincing view-dependent effects using directional dependence

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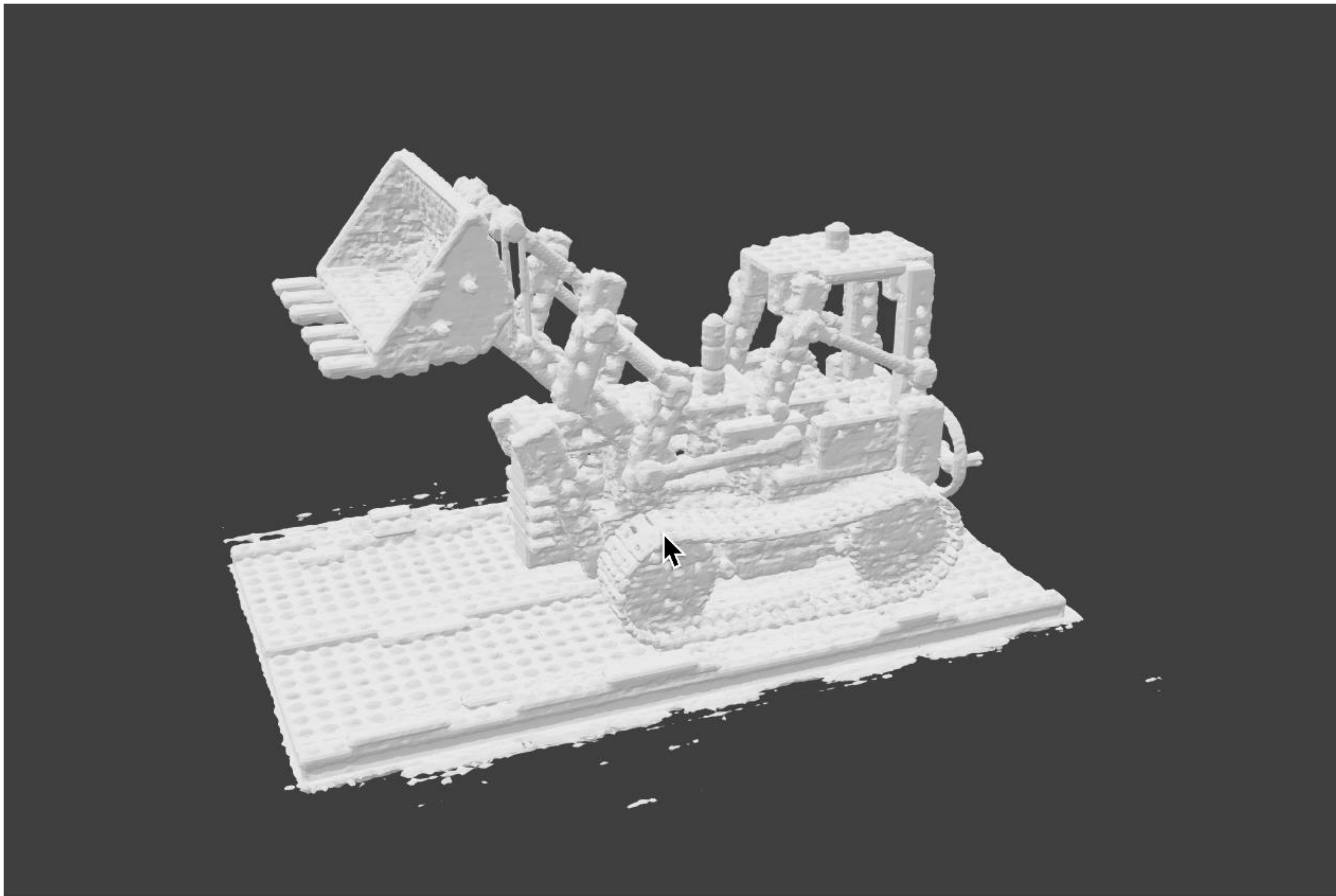
# NeRF encodes detailed scene geometry with occlusion effects

# NeRF encodes detailed scene geometry with occlusion effects

# NeRF encodes detailed scene geometry with occlusion effects

# NeRF encodes detailed scene geometry

# NeRF: 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)

▸ Apply positional encoding before passing coordinates into network to recover high frequency details