# Scene Representation Networks:
## Continuous 3D-Structure-Aware Neural Scene Representations

Vincent Sitzmann    Michael Zollhöfer    Gordon Wetzstein
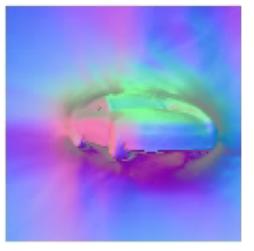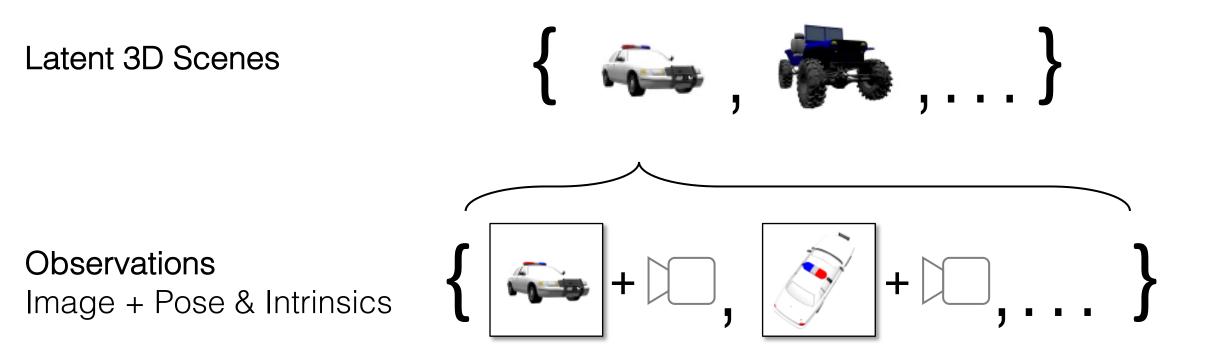
single image
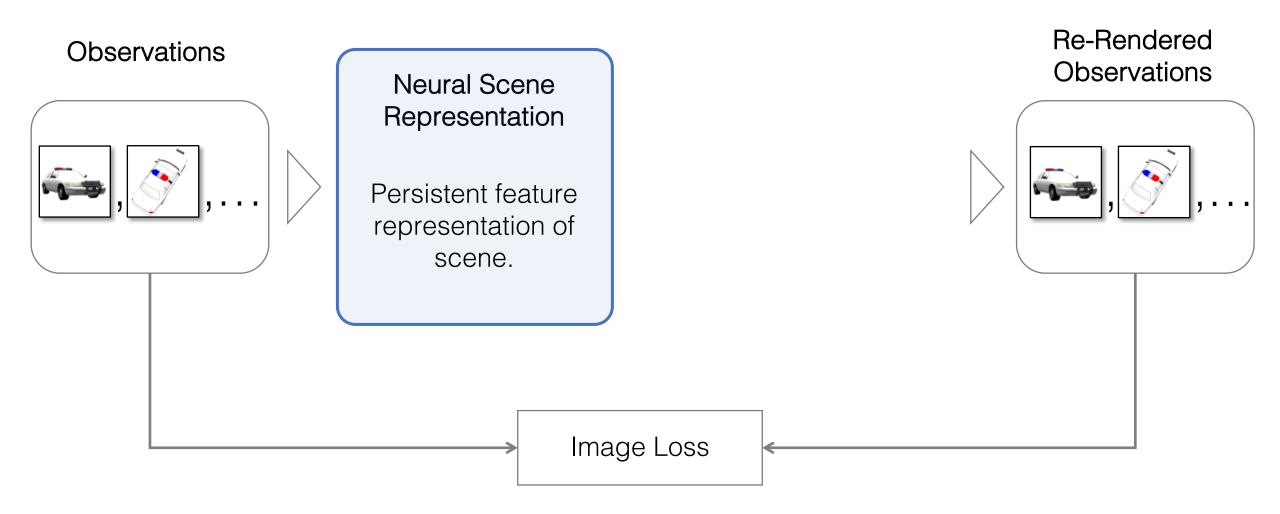camera pose
intrinsics

Novel Views          Surface Normals
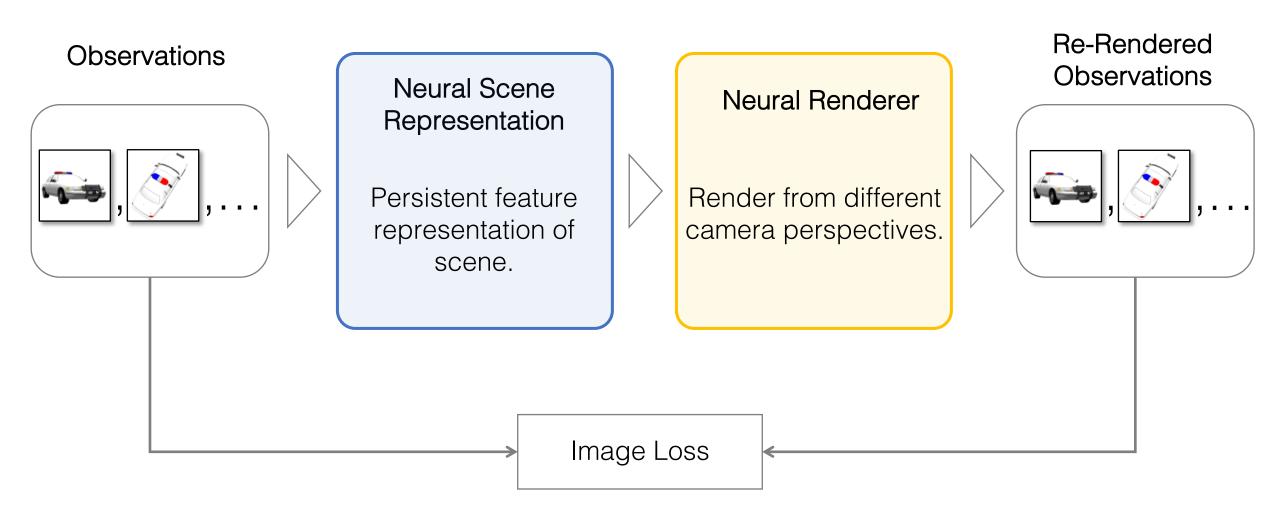
# Self-supervised Scene Representation Learning

**Latent 3D Scenes**

{  ,  , . . . }

**Observations**
Image + Pose & Intrinsics

{  +  ,  +  , . . . }

What can we learn about latent 3D scenes from observations?

Vision: Learn rich representations just by watching video!

# Self-supervised Scene Representation Learning

# Self-supervised Scene Representation Learning

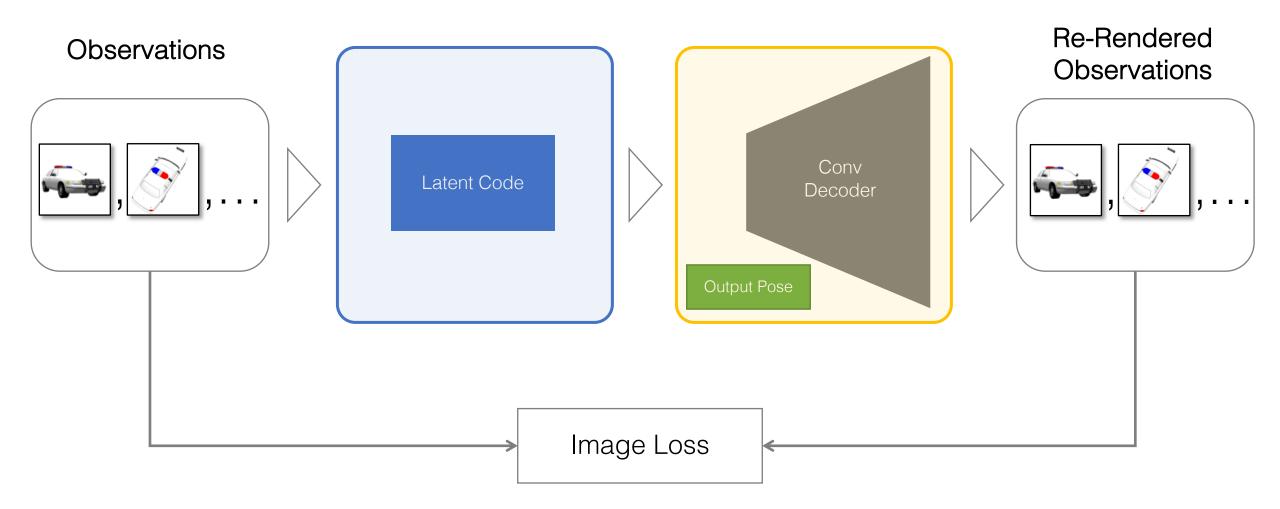# Self-supervised Scene Representation Learning

# 2D baseline: Autoencoder

# 2D baseline: Autoencoder

Doesn't capture 3D properties of scenes.

Trained on ~2500 shapenet cars with 50 observations each.
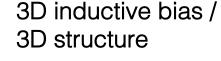


Need 3D inductive bias!

# Related Work



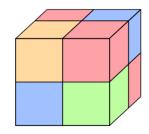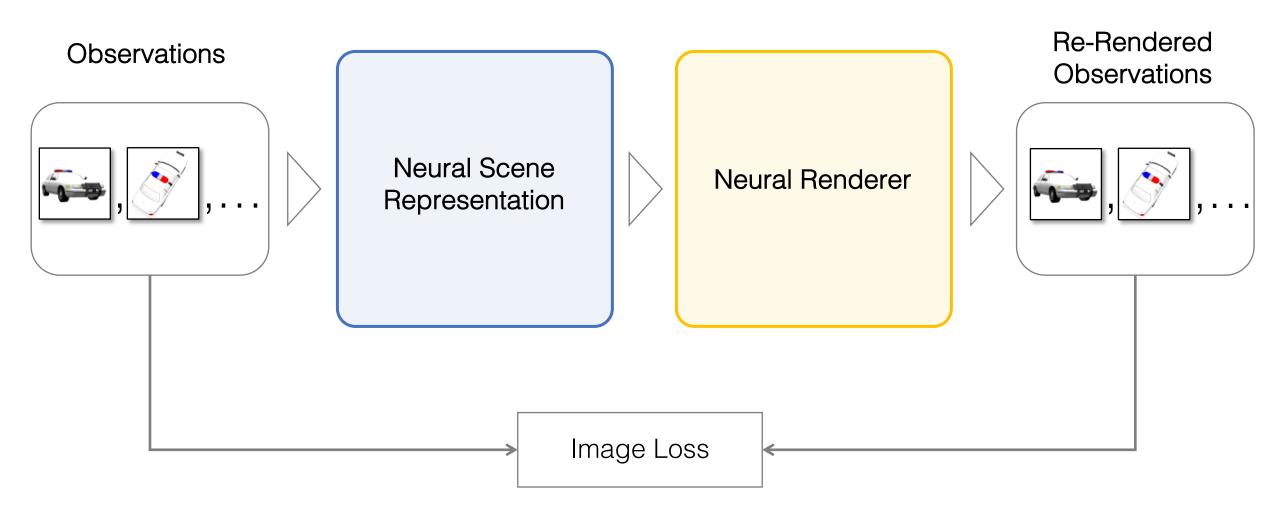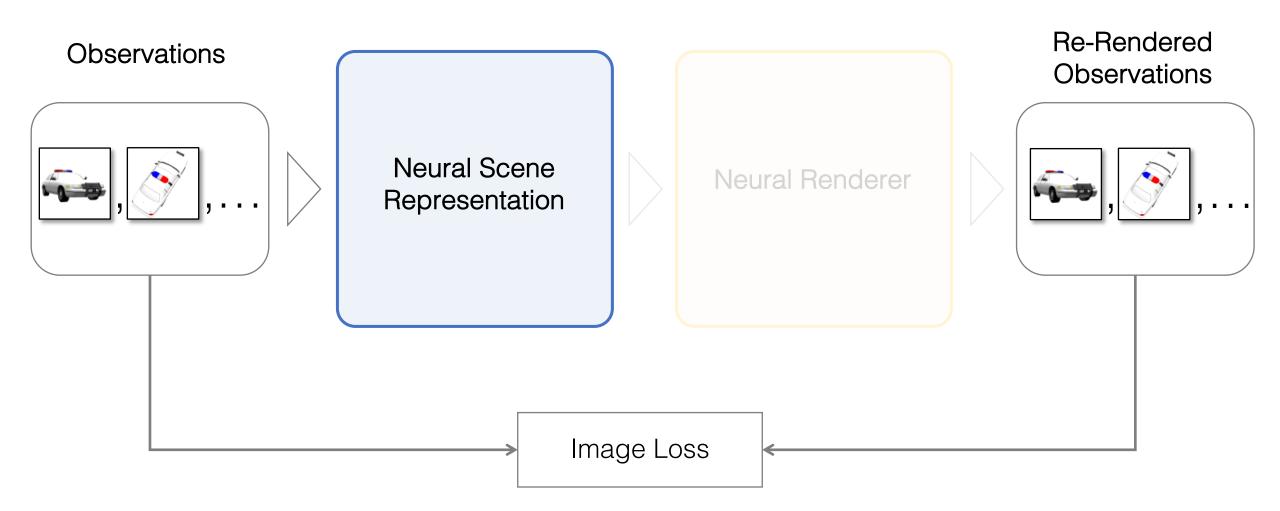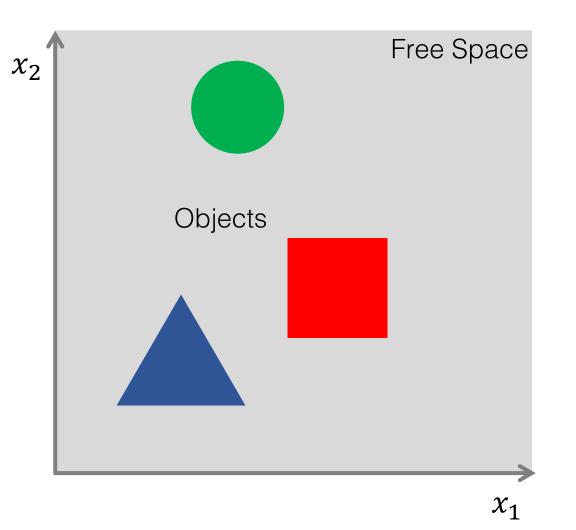|  | 3D inductive bias / 3D structure | Self-supervised with posed images |
|---|---|---|
| **Scene Representation Learning**<br>Tatarchenko et al., 2015<br>Worrall et al., 2017<br>Eslami et al., 2018<br>… | ✗ | ✓ |
| **2D Generative Models**<br>Goodfellow et al., 2014<br>Kingma et al., 2013<br>Kingma et al., 2018<br>… | ✗ | ✓ |
| **3D Computer Vision**<br>Choy et al., 2016<br>Huang et al., 2018<br>Park et al., 2018<br>… | ✓ | ✗ |
| **Voxel-based Representations**<br>Sitzmann et al., 2019<br>Lombardi et al., 2019<br>Phuoc et al., 2019<br>… |  |  |

- Memory inefficient: $O(n^3)$.
- Doesn't parameterize scene surfaces smoothly.
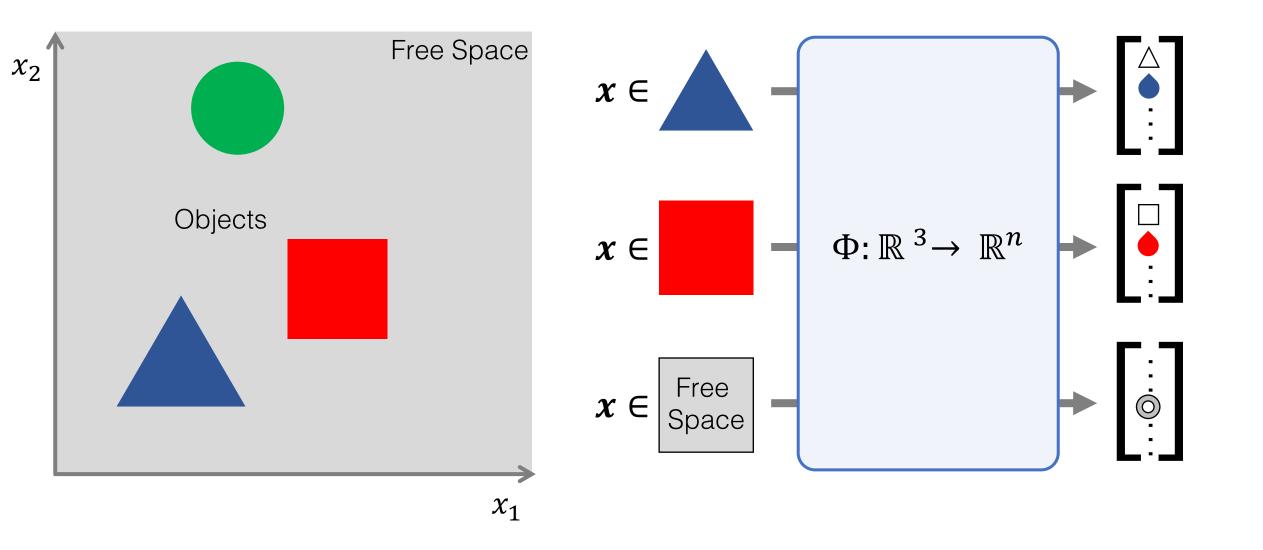- Generalization is hard.

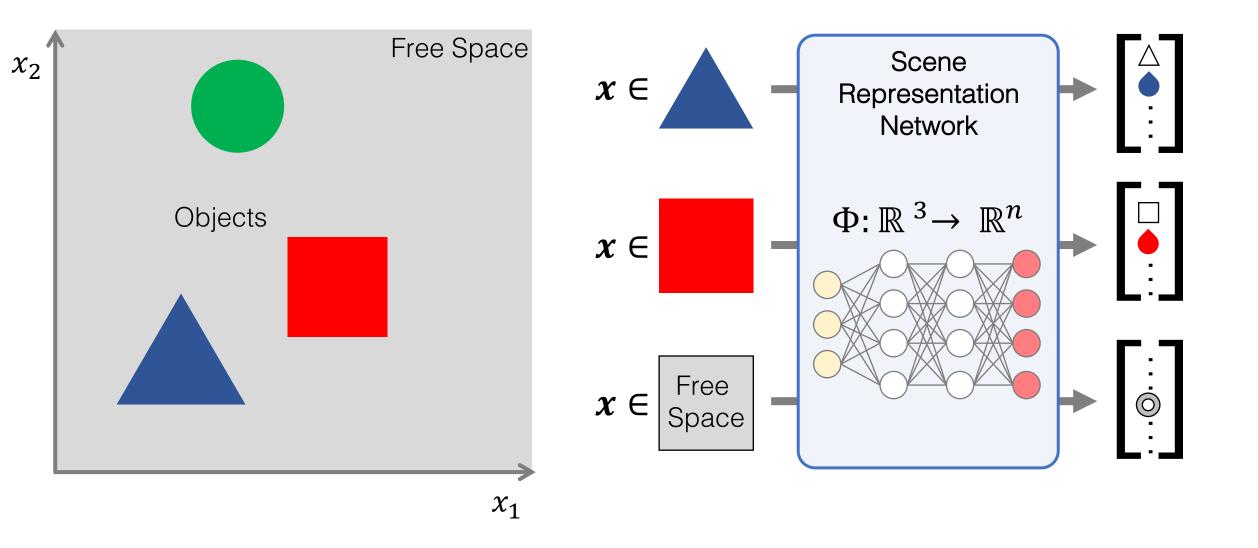# Scene Representation Networks

# Scene Representation Networks

# Model scene as function Φ that maps coordinates to features.

# Scene Representation Network parameterizes Φ as MLP.
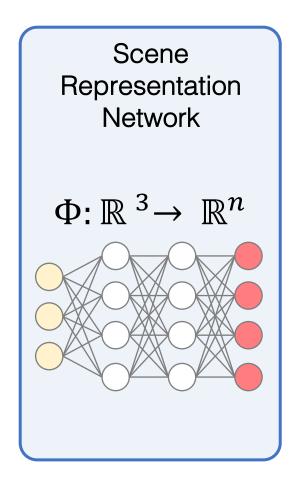
# Scene Representation Network parameterizes Φ as MLP.
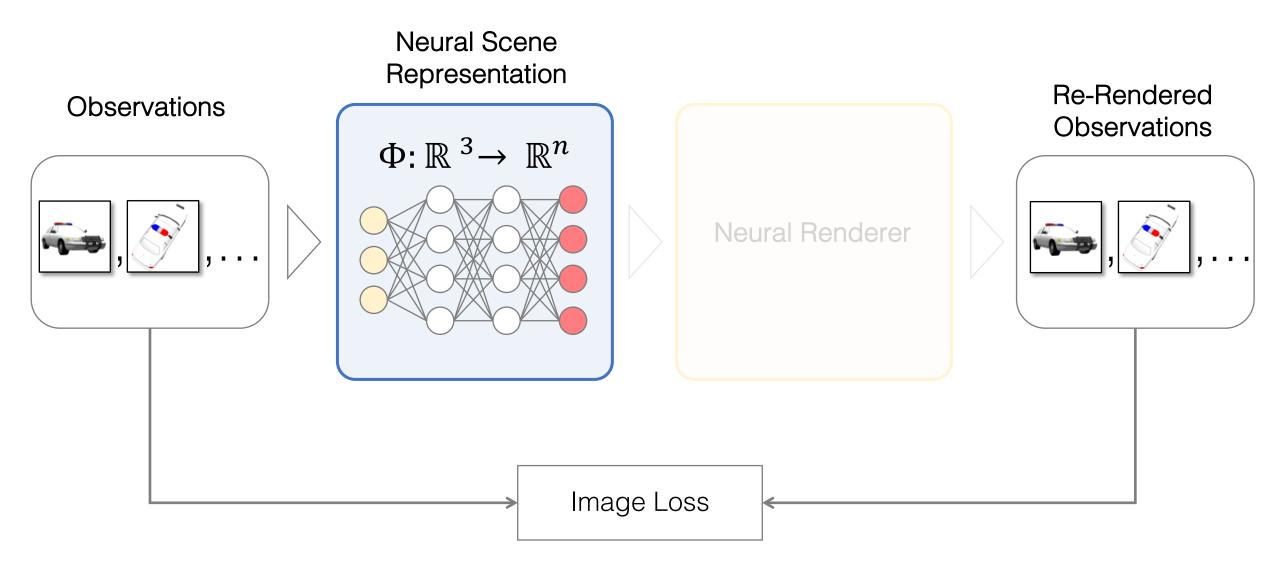
Can sample anywhere,
at arbitrary resolutions.

Parameterizes scene
surfaces smoothly.

Memory scales with scene
complexity.

Scene
Representation
Network

$$\Phi: \mathbb{R}^3 \to \mathbb{R}^n$$

# Scene Representation Networks

Observations

Neural Scene Representation

$$\Phi: \mathbb{R}^3 \rightarrow \mathbb{R}^n$$

Neural Renderer

Re-Rendered Observations

Image Loss

# Scene Representation Networks



Observations

Neural Scene Representation

$\Phi : \mathbb{R}^3 \to \mathbb{R}^n$

Neural Renderer

Re-Rendered Observations

Image Loss

# Neural Renderer.
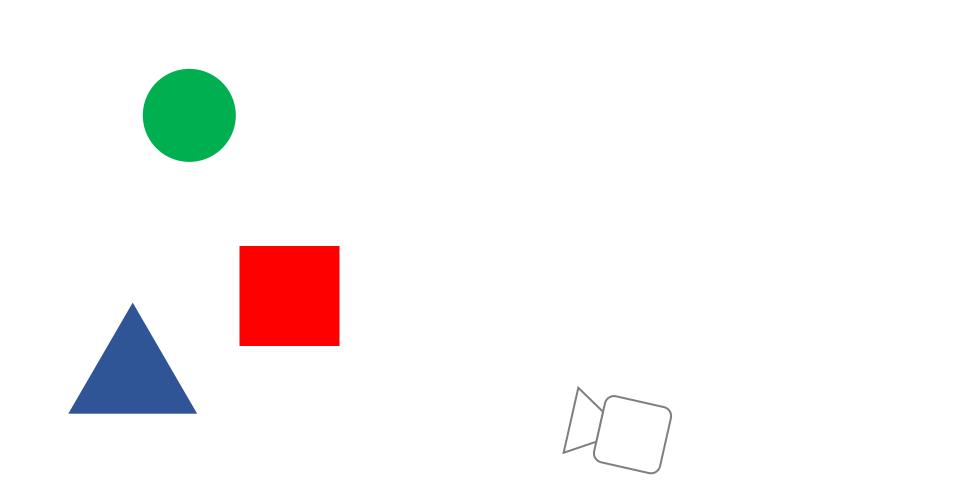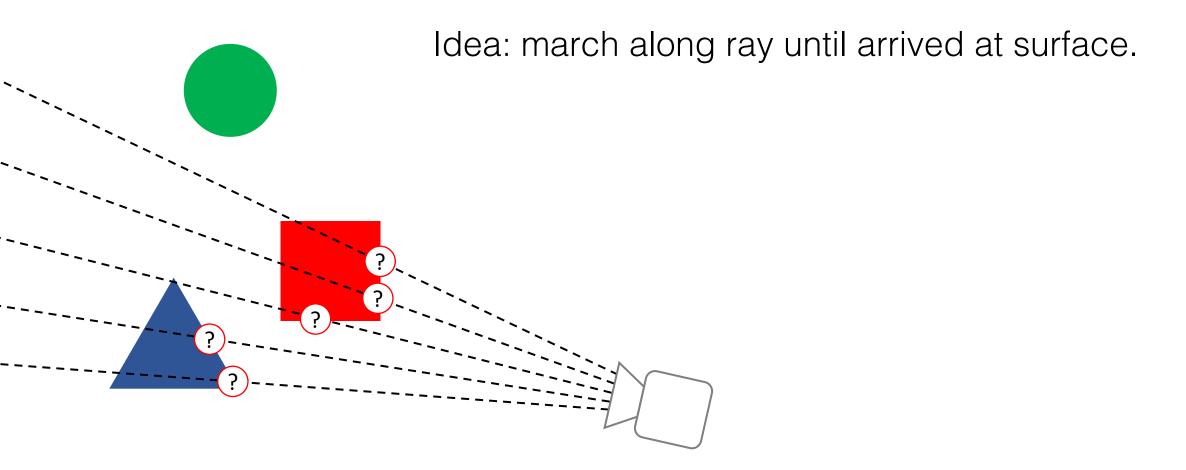
# Neural Renderer.

# Neural Renderer.

# Neural Renderer Step 1: Intersection Testing.

Idea: march along ray until arrived at surface.

# Neural Renderer Step 1: Intersection Testing.

# Neural Renderer Step 1: Intersection Testing.



**Ray Marching LSTM**

$\delta_{i+1}$
Step length

$v_i$ feature vector

Scene Representation
$\Phi: \mathbb{R}^3 \rightarrow \mathbb{R}^n$

Feasible step length:
Distance to closest scene surface

$\mathbf{x}_{i+1}$

$\mathbf{x_0}$

$\mathbf{x}_i$
world coordinates

# Neural Renderer Step 1: Intersection Testing.

Iteration 0

# Neural Renderer Step 1: Intersection Testing.

Iteration 1

# Neural Renderer Step 1: Intersection Testing.

Iteration 2

# Neural Renderer Step 1: Intersection Testing.

Iteration 3

# Neural Renderer Step 2: Color Generation

Iteration 4

# Neural Renderer Step 1: Intersection Testing.

Iteration …

# Neural Renderer Step 1: Intersection Testing.

# Neural Renderer Step 2: Color Generation



Scene Representation
$$\Phi: \mathbb{R}^3 \to \mathbb{R}^n$$

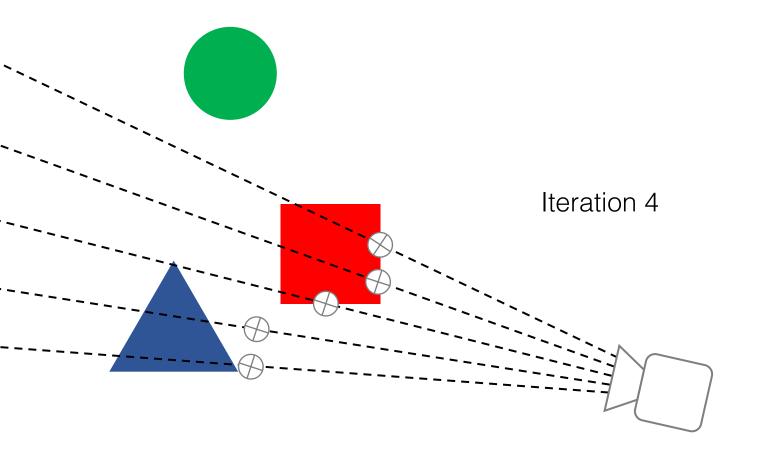Color MLP

# Can now train end-to-end with posed images only!

Observations

Neural Scene Representation

$$\Phi: \mathbb{R}^3 \to \mathbb{R}^n$$

Neural Renderer

Re-Rendered Observations

Image Loss

Generalizing across a class of scenes

# Each scene represented by its own SRN.

parameters $\phi_0 \in \mathbb{R}^l$

parameters $\phi_1 \in \mathbb{R}^l$

parameters $\phi_2 \in \mathbb{R}^l$

$\circ\ \circ\ \circ$

parameters $\phi_n \in \mathbb{R}^l$

# Each scene represented by its own SRN.

parameters $\phi_0 \in \mathbb{R}^l$

parameters $\phi_1 \in \mathbb{R}^l$

$\phi_i$ live on k-dimensional subspace of $\mathbb{R}^l$, $k < l$.

parameters $\phi_2 \in \mathbb{R}^l$

○ ○ ○

parameters $\phi_n \in \mathbb{R}^l$

# Each scene represented by its own SRN.

embedding $z_0 \in \mathbb{R}^k$

embedding $z_1 \in \mathbb{R}^k$

Represent each scene with
low-dimensional embedding

embedding $z_2 \in \mathbb{R}^k$

○ ○ ○

embedding $z_n \in \mathbb{R}^k$

parameters $\phi_0 \in \mathbb{R}^l$

parameters $\phi_1 \in \mathbb{R}^l$

parameters $\phi_2 \in \mathbb{R}^l$

○ ○ ○

parameters $\phi_n \in \mathbb{R}^l$

# Each scene represented by its own SRN.



embedding $z_0 \in \mathbb{R}^k$

embedding $z_1 \in \mathbb{R}^k$

embedding $z_2 \in \mathbb{R}^k$

embedding $z_n \in \mathbb{R}^k$

Hypernetwork

$$\Psi : \mathbb{R}^k \to \mathbb{R}^l,$$
$$z_i \mapsto \Psi(z_i) = \phi_i$$

parameters $\phi_0 \in \mathbb{R}^l$

parameters $\phi_1 \in \mathbb{R}^l$

parameters $\phi_2 \in \mathbb{R}^l$

parameters $\phi_n \in \mathbb{R}^l$

# Results

# Novel View Synthesis – Baseline Comparison

Shapenet v2 – *single-shot reconstruction* of objects in held-out test set

## Training

- Shapenet cars / chairs.

- 50 observations per object.

## Testing

- Cars / chairs from unseen test set

- Single observation!

## Input pose

# Novel View Synthesis – SRN Output

Shapenet v2 – _single-shot reconstruction_ of objects in held-out test set
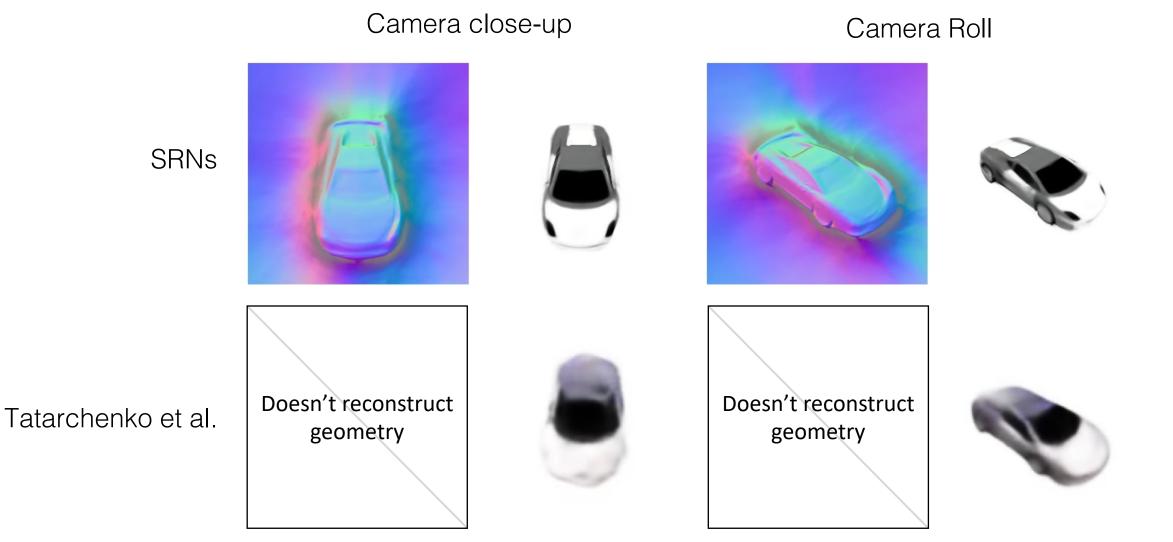
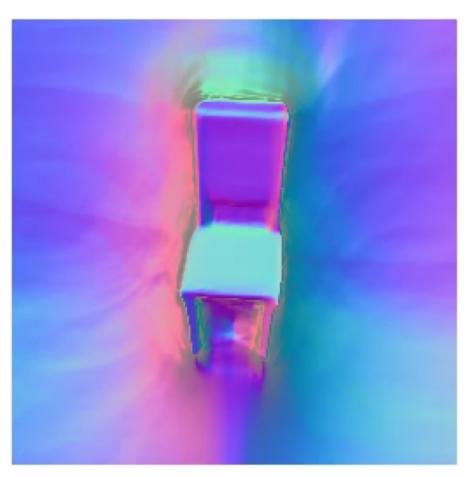# Sampling at arbitrary resolutions



32x32

64x64

128x128

256x256

512x512

Surface Normals

RGB

# Generalization to unseen camera poses

Camera close-up

Camera Roll

SRNs

# Generalization to unseen camera poses

# Latent code interpolation



Surface Normals

RGB

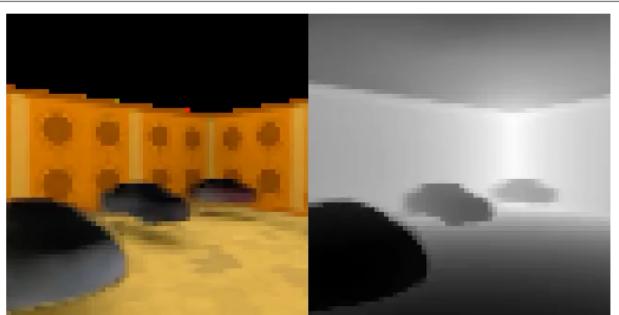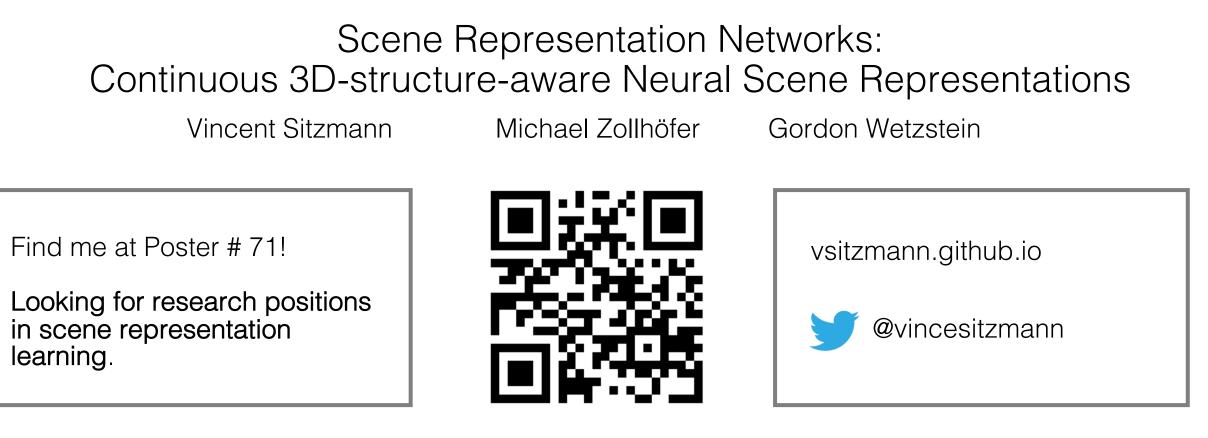# Latent code interpolation



Surface Normals

RGB

# Can represent room-scale scenes, but aren't compositional.
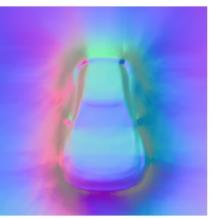


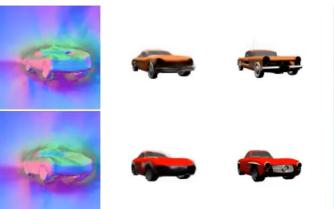Training set novel-view synthesis on GQN rooms (Eslami et al. 2018) with Shapenet cars, 50 observations.

Work-in-progress: Compositional SRNs generalize to unseen numbers of objects!