

# Implicit Neural Scene Representations and 3D-Aware Generative Modelling

Michael Niemeyer

Autonomous Vision Group  
MPI for Intelligent Systems and University of Tübingen

# Implicit Scene Representations for 3D Reconstruction

# Traditional 3D Reconstruction Pipeline

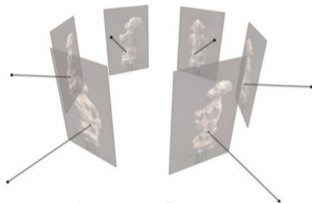


Input Images

# Traditional 3D Reconstruction Pipeline



Input Images

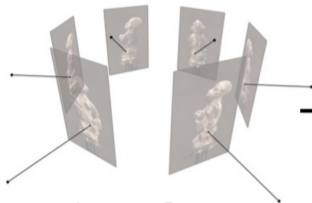


Camera Poses

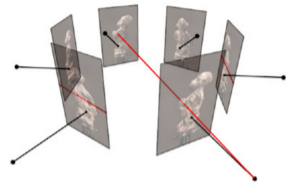
# Traditional 3D Reconstruction Pipeline



Input Images



Camera Poses

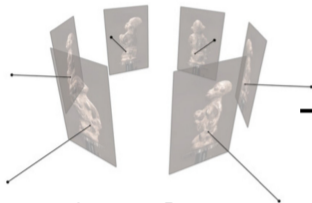


Dense Correspondences

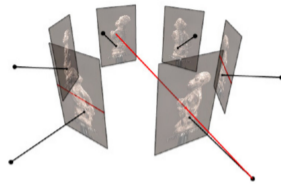
# Traditional 3D Reconstruction Pipeline



Input Images



Camera Poses



Dense Correspondences

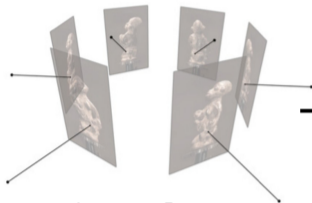


Depth Maps

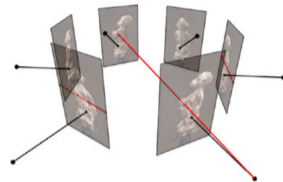
# Traditional 3D Reconstruction Pipeline



Input Images



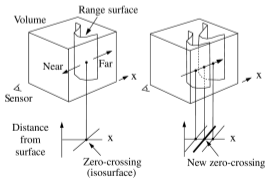
Camera Poses



Dense Correspondences

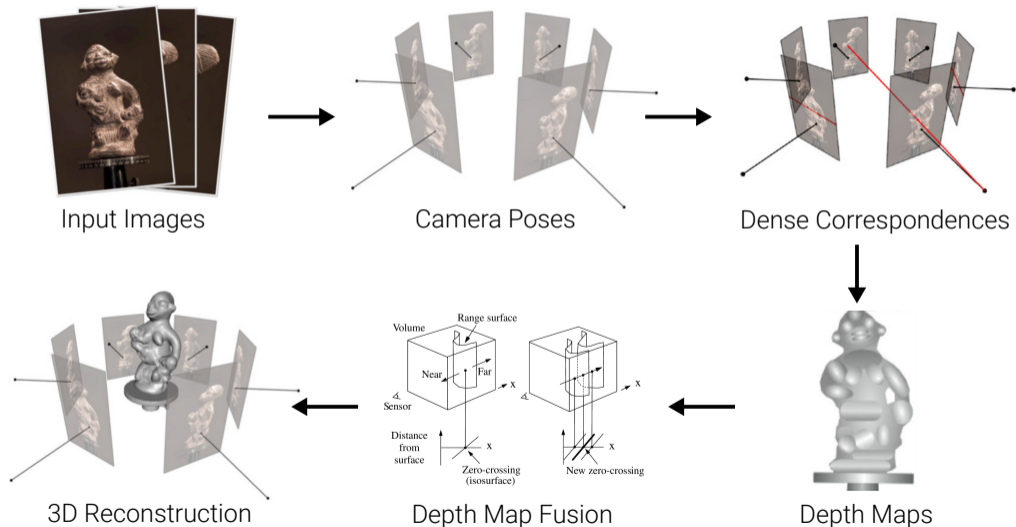


Depth Maps



Depth Map Fusion

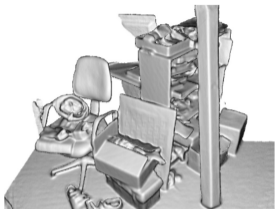
# Traditional 3D Reconstruction Pipeline



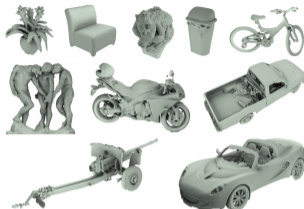


Can we **learn** 3D reconstruction **from data**?

# 3D Datasets and Repositories



[Newcombe et al., 2011]



[Choi et al., 2011]



[Dai et al., 2017]



[Wu et al., 2015]



[Chang et al., 2015]

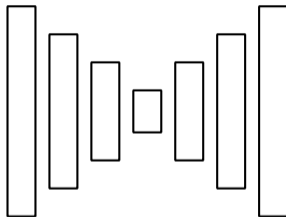


[Chang et al., 2017]

# 3D Reconstruction from a 2D Image



Input Images



Neural Network



3D Reconstruction

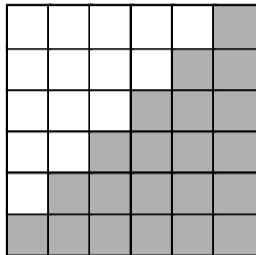
What is a good output representation?

# Learning-based 3D Reconstruction

## Voxels:

- ▶ **Discretization** of 3D space into grid
- ▶ Easy to process with neural networks
- ▶ Cubic memory  $O(n^3) \Rightarrow$  limited resolution
- ▶ Manhattan world bias

[Maturana et al., IROS 2015]

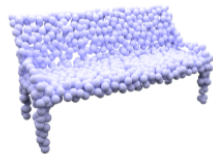
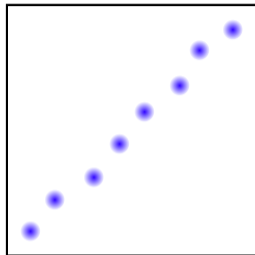


# Learning-based 3D Reconstruction

## Points:

- ▶ **Discretization** of surface into 3D points
- ▶ Does not model connectivity / topology
- ▶ Limited number of points
- ▶ Global shape description

[Fan et al., CVPR 2017]

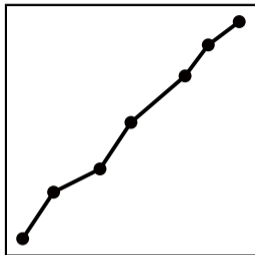


# Learning-based 3D Reconstruction

## Meshes:

- ▶ **Discretization** into vertices and faces
- ▶ Limited number of vertices / granularity
- ▶ Requires class-specific template – or –
- ▶ Leads to self-intersections

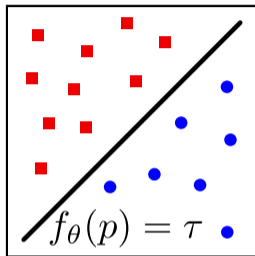
[Groueix et al., CVPR 2018]



# Learning-based 3D Reconstruction

## This work:

- ▶ Implicit representation  $\Rightarrow$  **No discretization**
- ▶ Arbitrary topology & resolution
- ▶ Low memory footprint
- ▶ Not restricted to specific class





# Occupancy Networks

## **Key Idea:**

- ▶ Do not represent 3D shape explicitly

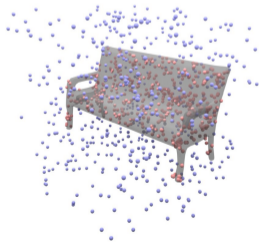
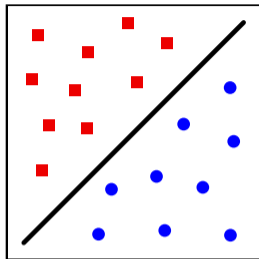
# Occupancy Networks

## Key Idea:

- ▶ Do not represent 3D shape explicitly
- ▶ Instead, consider surface **implicitly** as **decision boundary** of a non-linear classifier:

$$f_{\theta} : \mathbb{R}^3 \times \mathcal{X} \rightarrow [0, 1]$$

↑                    ↑                    ↑  
3D                    Condition                    Occupancy  
Location            (eg, Image)            Probability



# Occupancy Networks

## Key Idea:

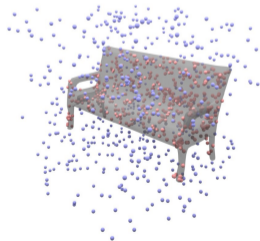
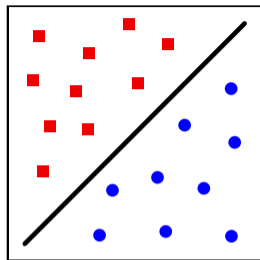
- ▶ Do not represent 3D shape explicitly
- ▶ Instead, consider surface **implicitly** as **decision boundary** of a non-linear classifier:

$$f_{\theta} : \mathbb{R}^3 \times \mathcal{X} \rightarrow [0, 1]$$

↑                    ↑                    ↑  
3D                    Condition                    Occupancy  
Location                    (eg, Image)                    Probability

## Remarks:

- ▶ The function  $f_{\theta}$  models an **occupancy field**



# Occupancy Networks

## Key Idea:

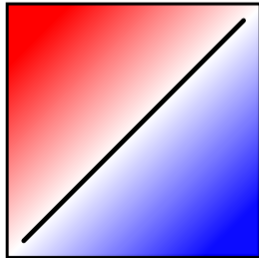
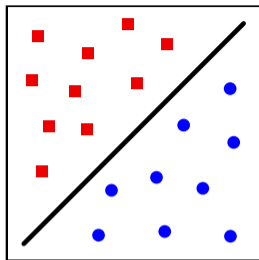
- ▶ Do not represent 3D shape explicitly
- ▶ Instead, consider surface **implicitly** as **decision boundary** of a non-linear classifier:

$$f_{\theta} : \mathbb{R}^3 \times \mathcal{X} \rightarrow [0, 1]$$

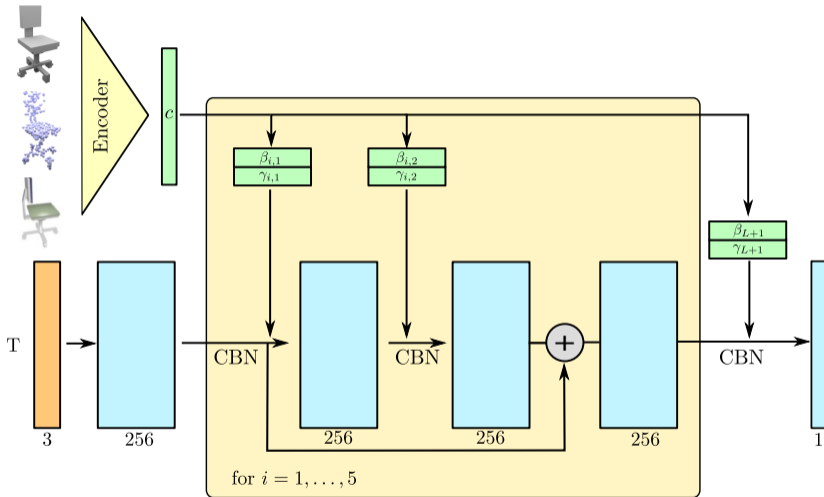
↑                    ↑                    ↑  
3D                    Condition                    Occupancy  
Location            (eg, Image)            Probability

## Remarks:

- ▶ The function  $f_{\theta}$  models an **occupancy field**
- ▶ Also possible: **signed distance field** [Park et al., 2019]



# Network Architecture



# Training Objective

## Occupancy Network:

$$\mathcal{L}(\theta, \psi) = \sum_{j=1}^K \text{BCE}(f_{\theta}(p_{ij}, z_i), o_{ij})$$

- ▶  $K$ : Randomly sampled 3D points ( $K = 2048$ )
- ▶ BCE: Cross-entropy loss

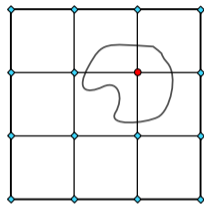
# Training Objective

## Variational Occupancy Encoder:

$$\mathcal{L}(\theta, \psi) = \sum_{j=1}^K \text{BCE}(f_{\theta}(p_{ij}, z_i), o_{ij}) + KL [q_{\psi}(z | (p_{ij}, o_{ij})_{j=1:K}) || p_0(z)]$$

- ▶  $K$ : Randomly sampled 3D points ( $K = 2048$ )
- ▶ BCE: Cross-entropy loss
- ▶  $q_{\psi}$ : Encoder

# Occupancy Networks



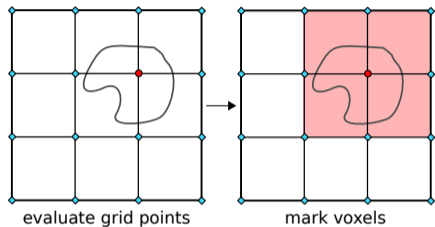
evaluate grid points

## **Multiresolution IsoSurface Extraction (MISE):**

- ▶ Build octree by incrementally querying the occupancy network



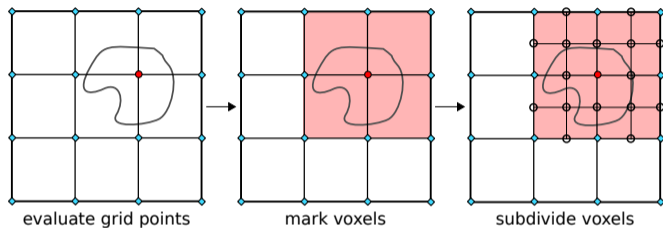
# Occupancy Networks



## Multiresolution IsoSurface Extraction (MISE):

- ▶ Build octree by incrementally querying the occupancy network

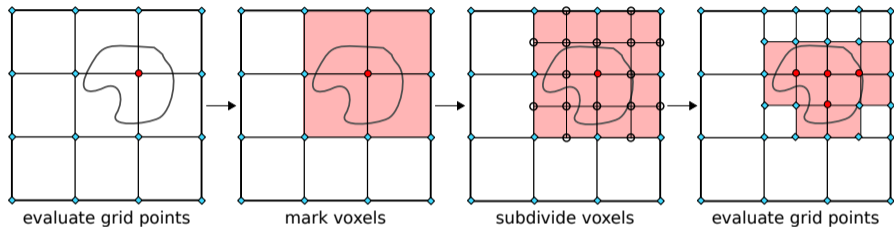
# Occupancy Networks



## Multiresolution IsoSurface Extraction (MISE):

- ▶ Build octree by incrementally querying the occupancy network

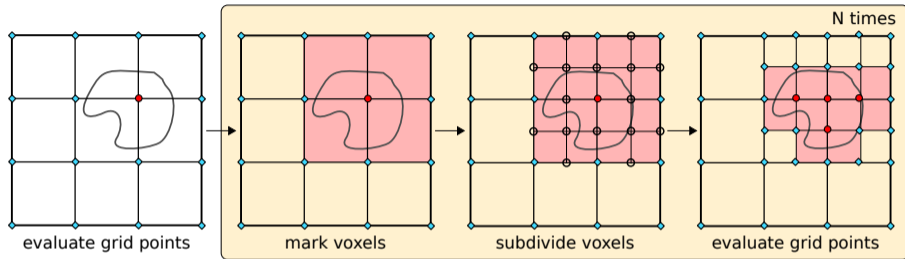
# Occupancy Networks



## Multiresolution IsoSurface Extraction (MISE):

- ▶ Build octree by incrementally querying the occupancy network

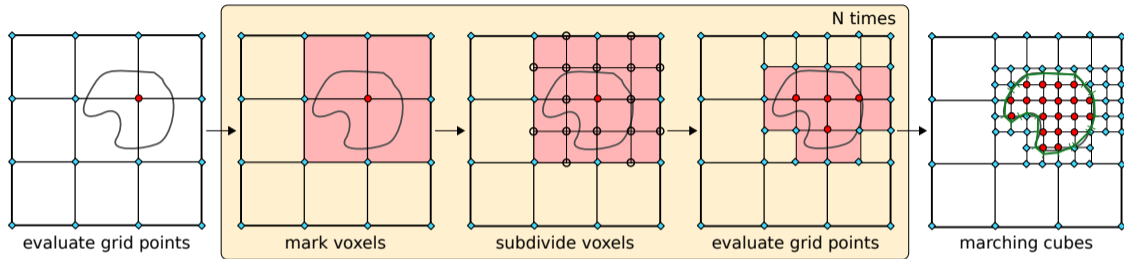
# Occupancy Networks



## Multiresolution IsoSurface Extraction (MISE):

- Build octree by incrementally querying the occupancy network

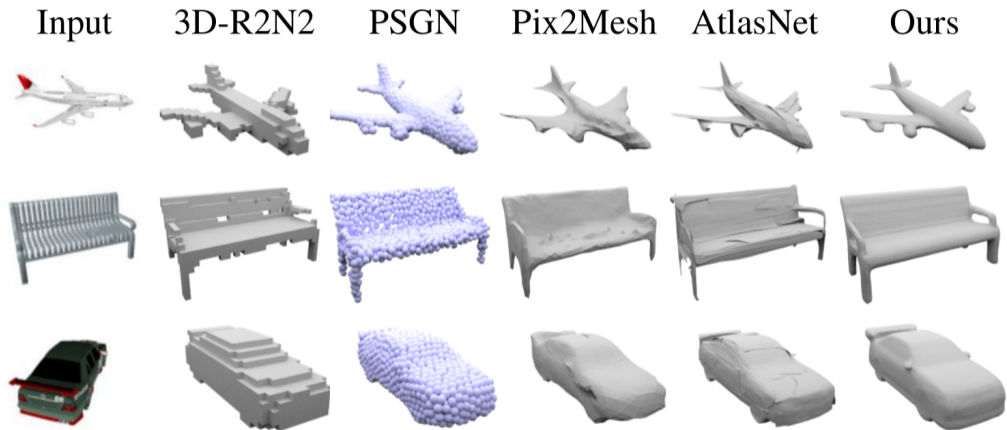
# Occupancy Networks



## Multiresolution IsoSurface Extraction (MISE):

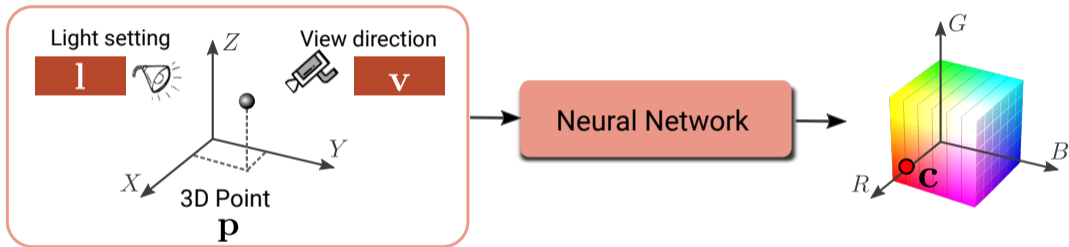
- ▶ Build octree by incrementally querying the occupancy network
- ▶ Extract triangular mesh using marching cubes algorithm (1-3 seconds in total)

# Results



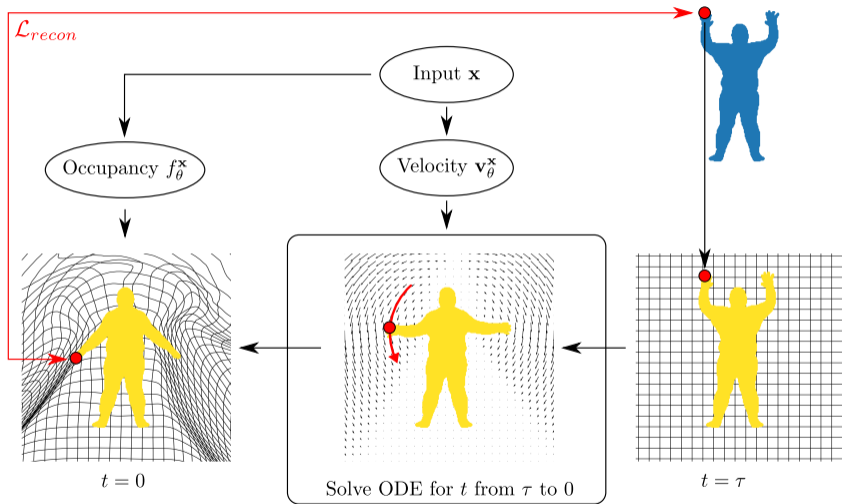
Applications

# Appearance



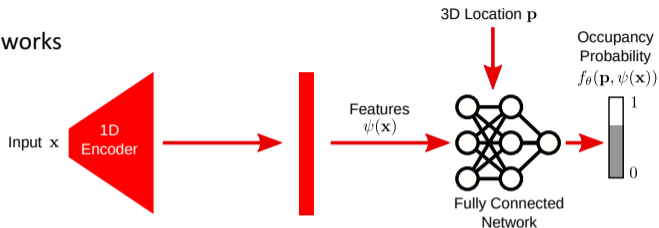


# Motion



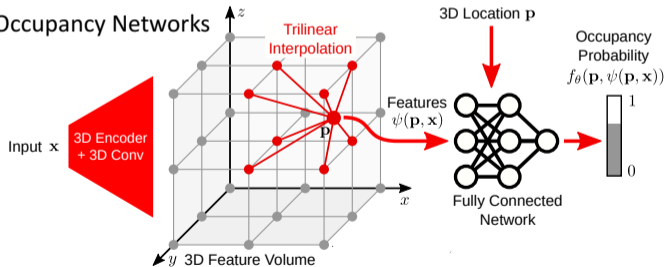
# 3D Scenes

## Occupancy Networks

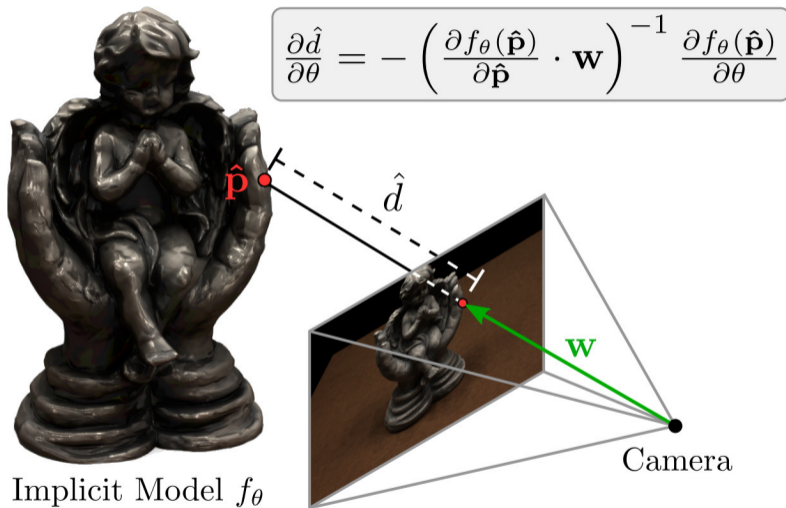


---

## Convolutional Occupancy Networks



# Differentiable Rendering



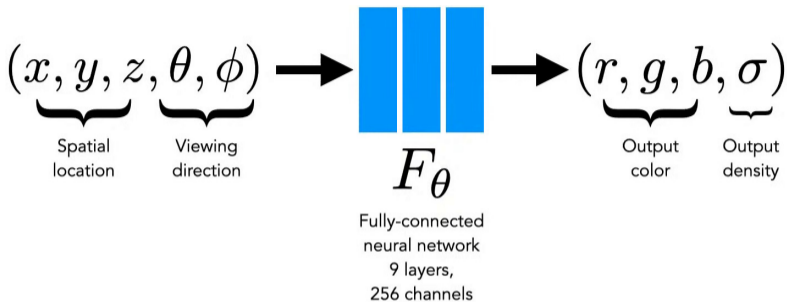
# Neural Rendering: Neural Radiance Fields

# Novel View Synthesis



- **Task:** Given a set of images of a scene (left), render image from novel viewpoint (right)

# NeRF: Representing Scenes as Neural Radiance Fields



- ▶ Vanilla **ReLU MLP** that maps from **location/view direction to color/density**
- ▶ **Density**  $\sigma$  describes how solid/transparent a 3D point is (can model, e.g., fog)
- ▶ Conditioning on view direction allows for modeling **view-dependent effects**

# Volume Rendering

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

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

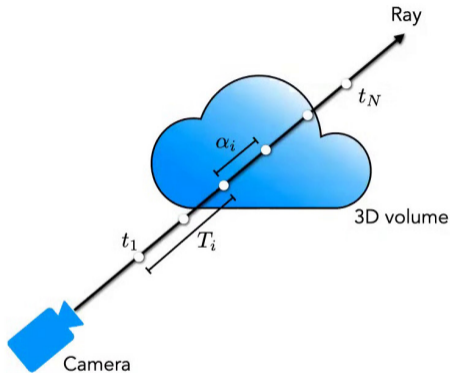
weights                      colors

How much light is blocked earlier along ray:

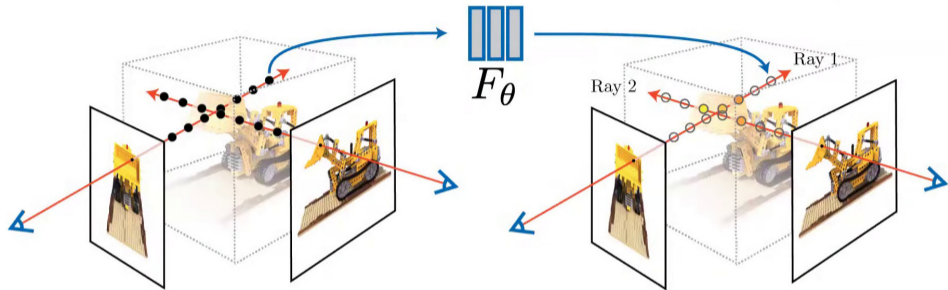
$$T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$$

How much light is contributed by ray segment  $i$ :

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



# NeRF Training



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

- Shoot ray, render ray to pixel, minimize **reconstruction error** via backpropagation



# Fourier Features



NeRF (Naive)



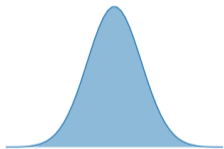
NeRF (with positional encoding)

- Essential trick: Compute **positional encoding** for input point  $\mathbf{x}$  and direction  $\mathbf{d}$

# Generative Neural Scene Representations

# Generative Models

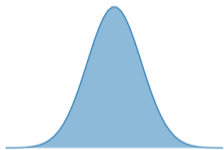
Sample a latent code from the prior distribution.



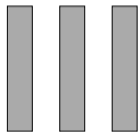
Latent Code

# Generative Models

Pass latent code to trained generator  $G_\theta$ .



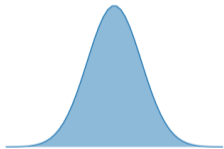
Latent Code



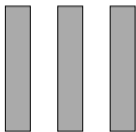
Generator  $G_\theta$

# Generative Models

The generator outputs a synthesized image.



Latent Code



Generator  $G_\theta$

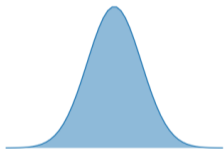


Generated Image\*

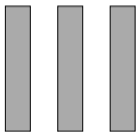
\* The generated images are samples from StyleGAN2.

# Generative Models

Sample more latent codes to get different generated images.



Latent Code



Generator  $G_\theta$

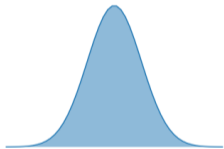


Generated Image\*

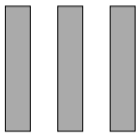
\* The generated images are samples from StyleGAN2.

# Generative Models

Sample more latent codes to get different generated images.



Latent Code



Generator  $G_\theta$



Generated Image\*

\* The generated images are samples from StyleGAN2.

Is the ability to sample photorealistic images  
all we want?



# Generative Models

For many applications, we require **control over the generation process**:

# Generative Models

For many applications, we require **control over the generation process**:

**Note:** This and the following videos are only shown when opened with a supported PDF reader (e.g. Okular).

## Animation Movies



Video Source: Disney's Toy Story 4 Trailer

# Generative Models

For many applications, we require **control over the generation process**:

## Video Games

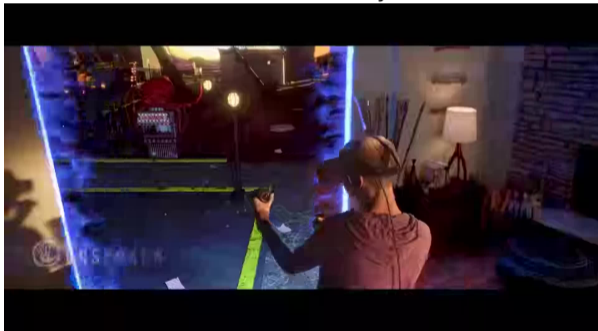


Video Source: Gran Turismo 7 Trailer

# Generative Models

For many applications, we require **control over the generation process**:

Virtual Reality



Video Source: Oculus Rift Trailer

# Generative Models

**Goal:** A generative model for **3D-aware image synthesis** which allows us to:

# Generative Models

**Goal:** A generative model for **3D-aware image synthesis** which allows us to:

- ▶ Generate photorealistic images

# Generative Models

**Goal:** A generative model for **3D-aware image synthesis** which allows us to:

- ▶ Generate photorealistic images
- ▶ Control individual objects wrt. their pose, size, and position in 3D

# Generative Models

**Goal:** A generative model for **3D-aware image synthesis** which allows us to:

- ▶ Generate photorealistic images
- ▶ Control individual objects wrt. their pose, size, and position in 3D
- ▶ Control camera viewpoint in 3D



# Generative Models

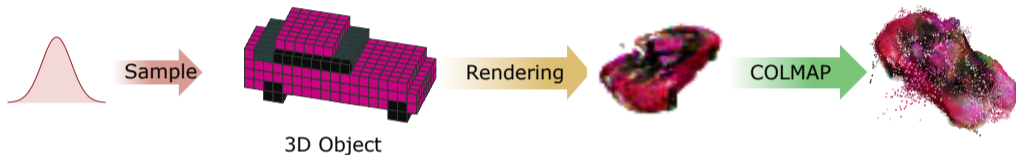
**Goal:** A generative model for **3D-aware image synthesis** which allows us to:

- ▶ Generate photorealistic images
- ▶ Control individual objects wrt. their pose, size, and position in 3D
- ▶ Control camera viewpoint in 3D
- ▶ Train from collections of unposed images

What representation should we use for  
3D-aware image synthesis?

# 3D Representations

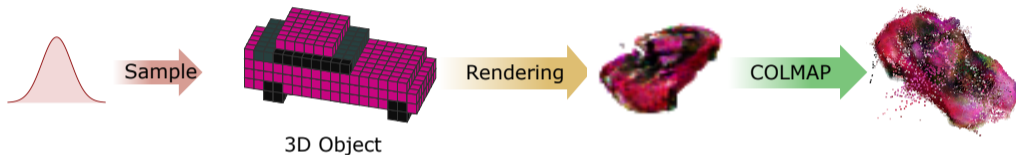
## Voxel-based 3D Shape with Volumetric Rendering



PlatonicGAN [Henzler et al., ICCV 2019]

# 3D Representations

## Voxel-based 3D Shape with Volumetric Rendering

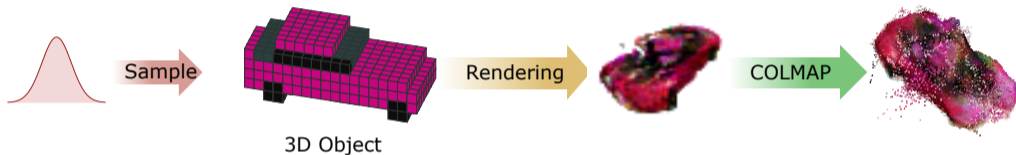


PlatonicGAN [Henzler et al., ICCV 2019]

+ Multi-view consistent

# 3D Representations

## Voxel-based 3D Shape with Volumetric Rendering

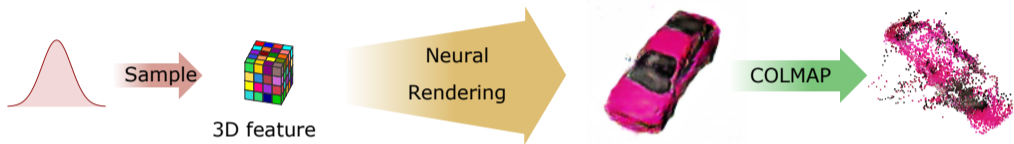


PlatonicGAN [Henzler et al., ICCV 2019]

- + Multi-view consistent
- Low image fidelity, high memory consumption

# 3D Representations

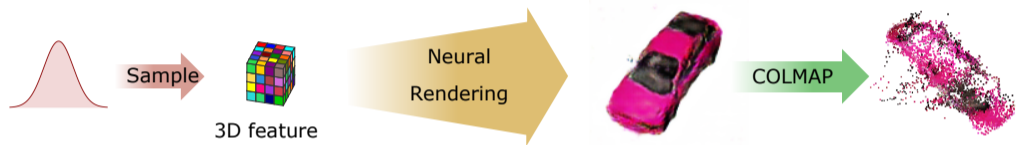
## Voxel-based 3D Latent Feature with Learnable Projection



HoloGAN [Nguyen-Phuoc et al., ICCV 2019]

# 3D Representations

## Voxel-based 3D Latent Feature with Learnable Projection

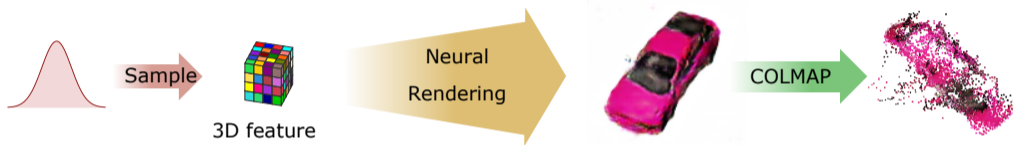


HoloGAN [Nguyen-Phuoc et al., ICCV 2019]

+ High image fidelity

# 3D Representations

## Voxel-based 3D Latent Feature with Learnable Projection



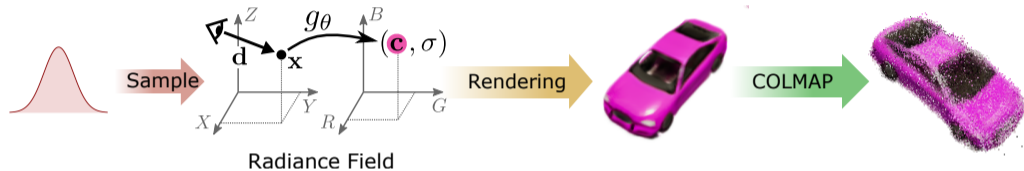
HoloGAN [Nguyen-Phuoc et al., ICCV 2019]

- + High image fidelity
- Object identity may vary with viewpoint due to learnable projection



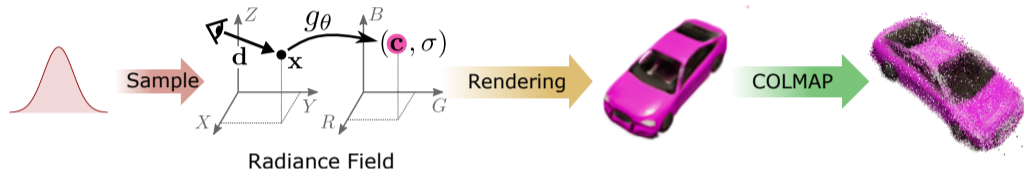
# 3D Representations

## Generative Radiance Fields



# 3D Representations

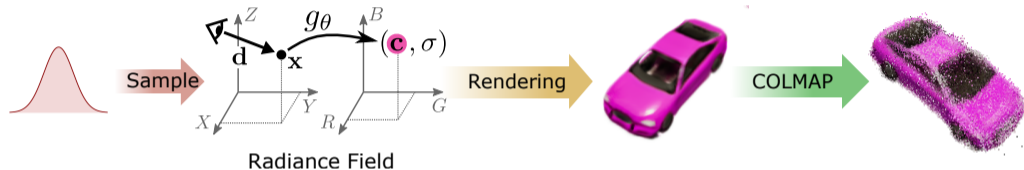
## Generative Radiance Fields



+ Continuous representation, multi-view consistent

# 3D Representations

## Generative Radiance Fields



- + Continuous representation, multi-view consistent
- + High image fidelity, low memory consumption

# Generative Radiance Fields

# Generative Radiance Fields

Sample camera matrix  $\mathbf{K}$ , camera pose  $\xi \sim p_\xi$ , and patch sampling pattern  $\nu \sim p_\nu$ .

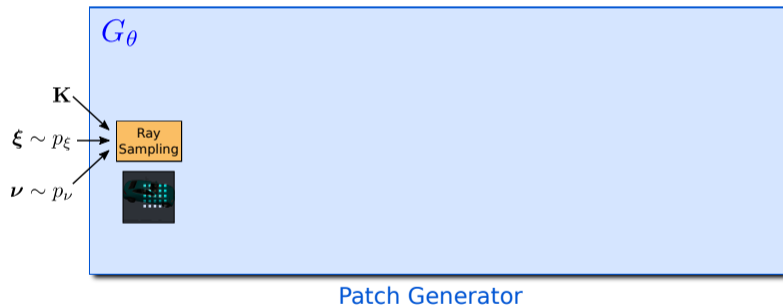
$\mathbf{K}$

$\xi \sim p_\xi$

$\nu \sim p_\nu$

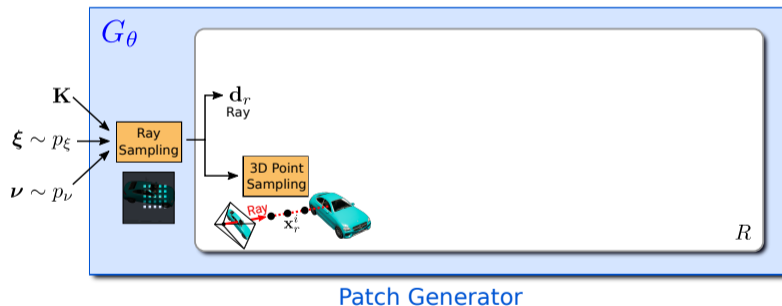
# Generative Radiance Fields

Pass  $\mathbf{K}$ ,  $\xi$ , and  $\nu$  to generator  $G_\theta$  and sample pixels / rays on image plane.



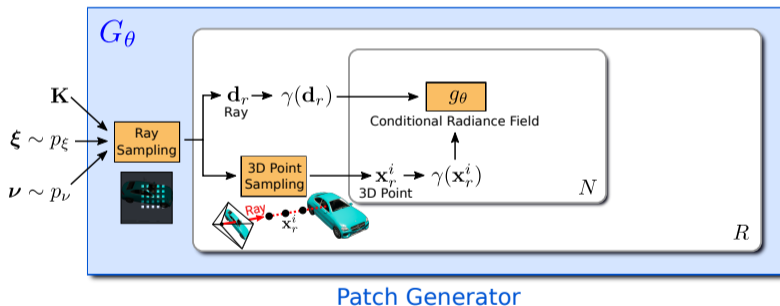
# Generative Radiance Fields

For each ray, get viewing direction  $\mathbf{d}_r$ , and sample 3D points  $\mathbf{x}_r^i$  along ray.



# Generative Radiance Fields

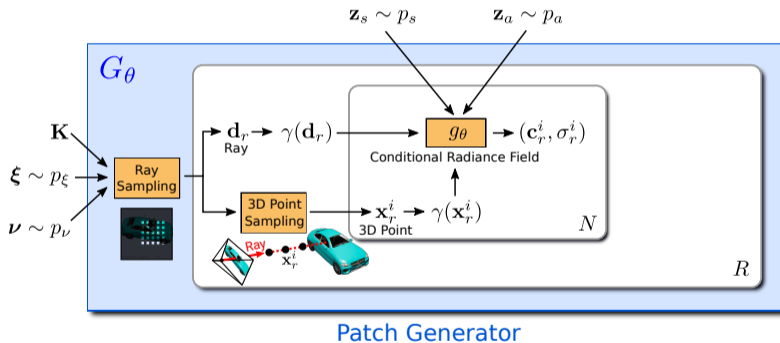
Pass  $\mathbf{d}_r$  and  $\mathbf{x}_r^i$  to positional encoding  $\gamma$  and then to the conditional radiance field  $g_\theta$ .





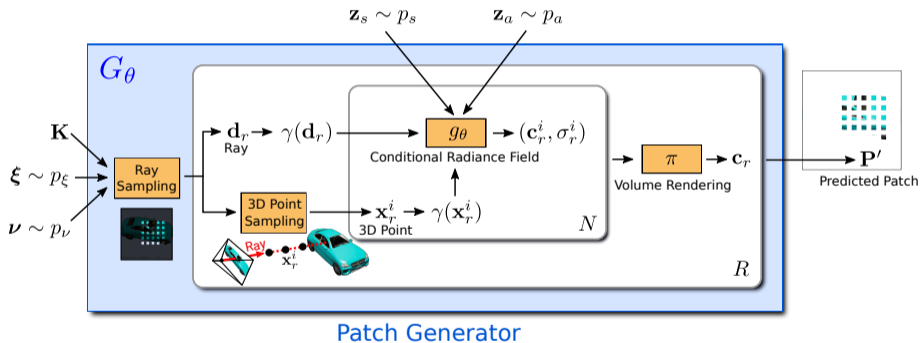
# Generative Radiance Fields

Sample latent shape and appearance codes  $\mathbf{z}_s, \mathbf{z}_a$  and pass them to  $g_\theta$ .



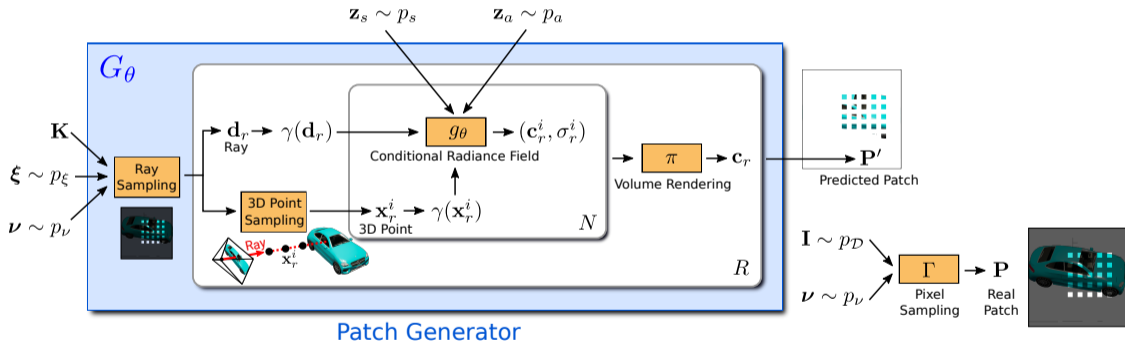
# Generative Radiance Fields

Perform volume-rendering for each ray and get predicted patch  $\mathbf{P}'$ .



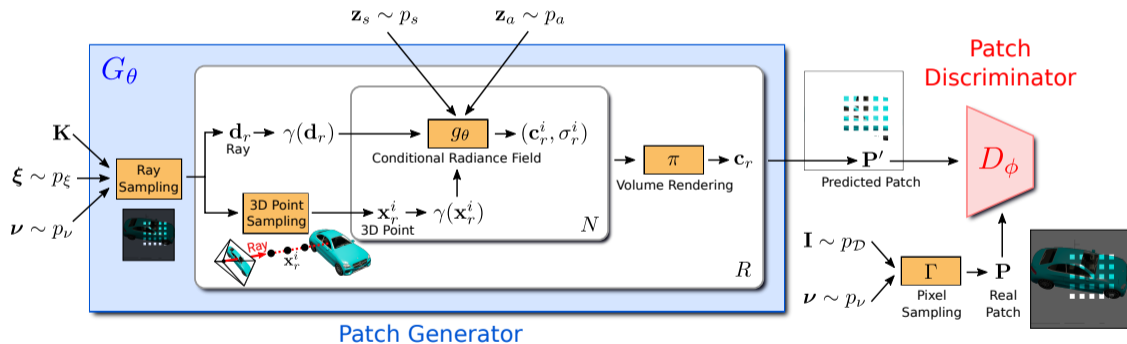
# Generative Radiance Fields

Sample patch  $\mathbf{P}$  from real image  $\mathbf{I}$  drawn from the data distribution  $p_{\mathcal{D}}$ .

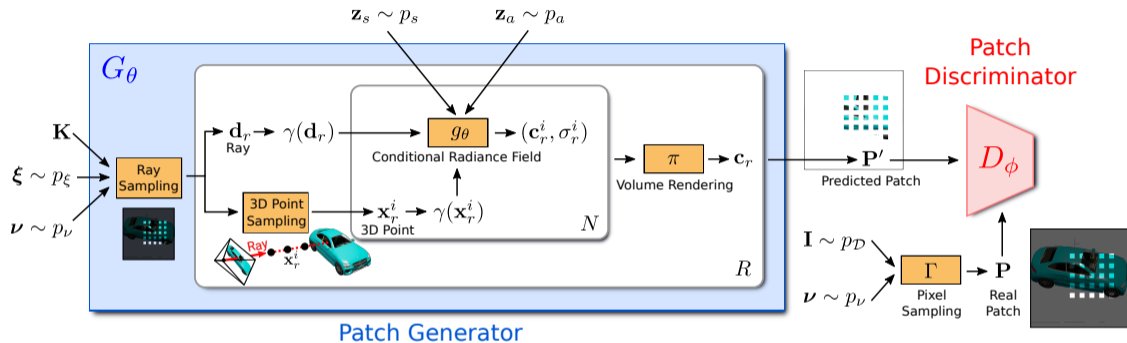


# Generative Radiance Fields

Pass fake and real patch  $\mathbf{P}'$ ,  $\mathbf{P}$  to discriminator  $D_\phi$  and train with adversarial loss.



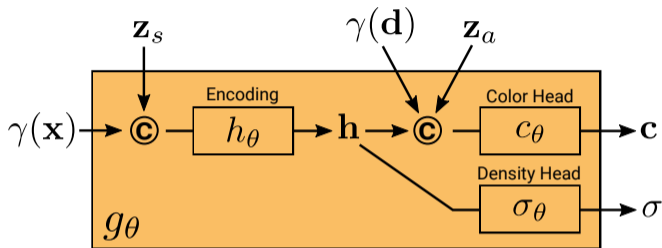
# Generative Radiance Fields



- ▶ Generator/discriminator for **image patches** of size  $32 \times 32$  pixels
- ▶ Patches sampled at **random scale** using dilation

How do we parametrize  
Conditional Radiance Fields?

# Conditional Radiance Fields



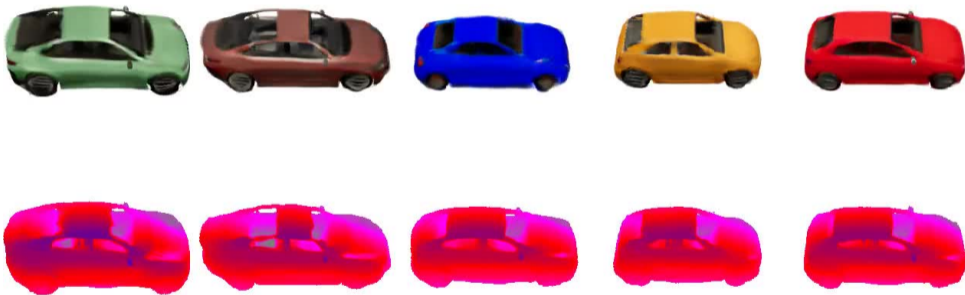
- ▶ Conditional radiance fields as fully-connected MLPs with ReLU activation
- ▶ Shape code  $\mathbf{z}_s$  concatenated with encoded 3D location  $\gamma(\mathbf{x})$
- ▶ Appearance code  $\mathbf{z}_a$  concatenated with encoded viewing direction  $\gamma(\mathbf{d})$

How well does it work?



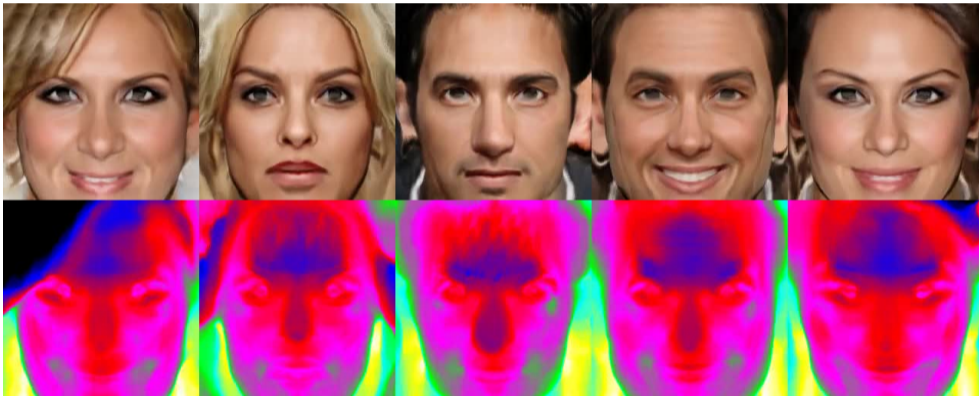
# Generative Radiance Fields

Results on synthetic Carla dataset at  $256^2$  pixels:



# Generative Radiance Fields

Results on real CelebA-HQ dataset at  $256^2$  pixels:



How can we scale to  
more complex, multi-object scenes?

# GIRAFFE: Compositional Generative Neural Feature Fields

GRAF:

- ▶ Incorporate a **3D representation** into the generative model

# GIRAFFE: Compositional Generative Neural Feature Fields

GRAF:

- ▶ Incorporate a **3D representation** into the generative model

GIRAFFE:

- ▶ Incorporate a **compositional 3D scene representation** into the generative model

# GIRAFFE: Compositional Generative Neural Feature Fields

GRAF:

- ▶ Incorporate a **3D representation** into the generative model

GIRAFFE:

- ▶ Incorporate a **compositional 3D scene representation** into the generative model
- ▶ Incorporate a **neural renderer** to yield fast and high-quality inference

GIRAFFE

# GIRAFFE

Sample  $N$  shape and appearance codes.



Shape and  
Appearance



Shape and  
Appearance

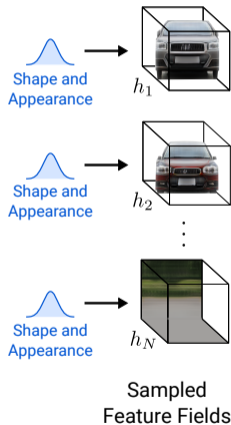


Shape and  
Appearance



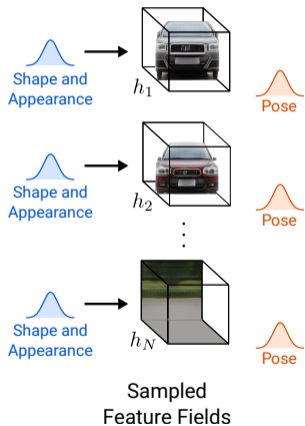
# GIRAFFE

Get  $N$  feature fields. Note: We show features in RGB color for clarity.



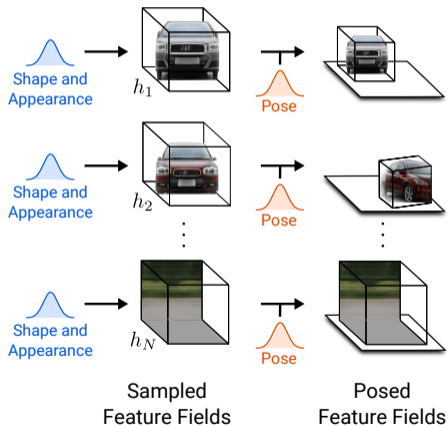
# GIRAFFE

Sample size and pose for each feature field.



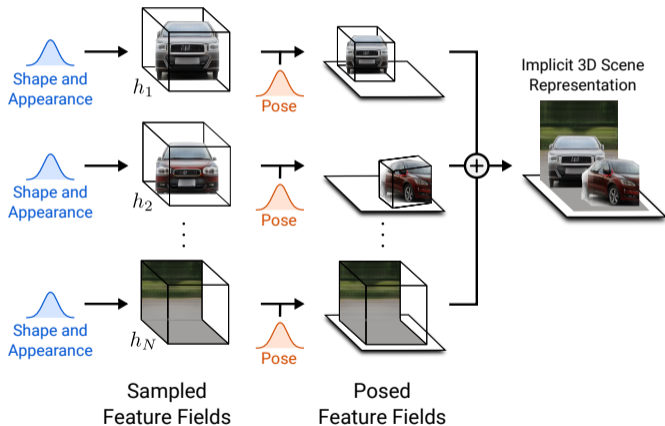
# GIRAFFE

Get posed feature fields.



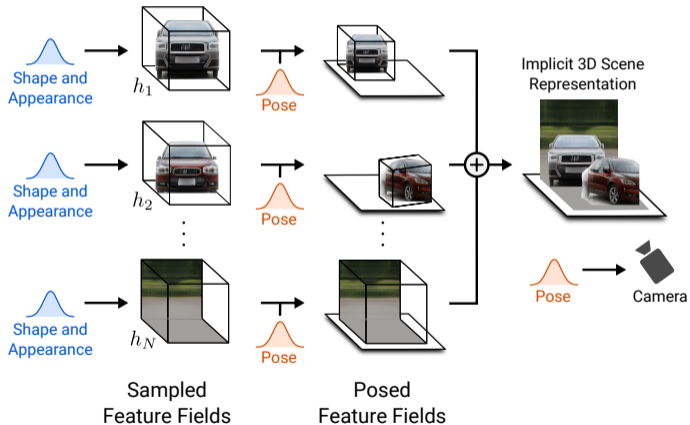
# GIRAFFE

Composite all feature feature fields to one 3D scene representation.



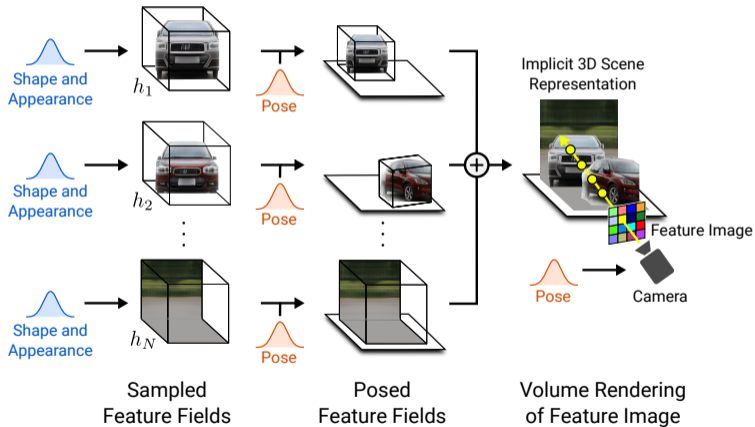
# GIRAFFE

Sample a camera pose.



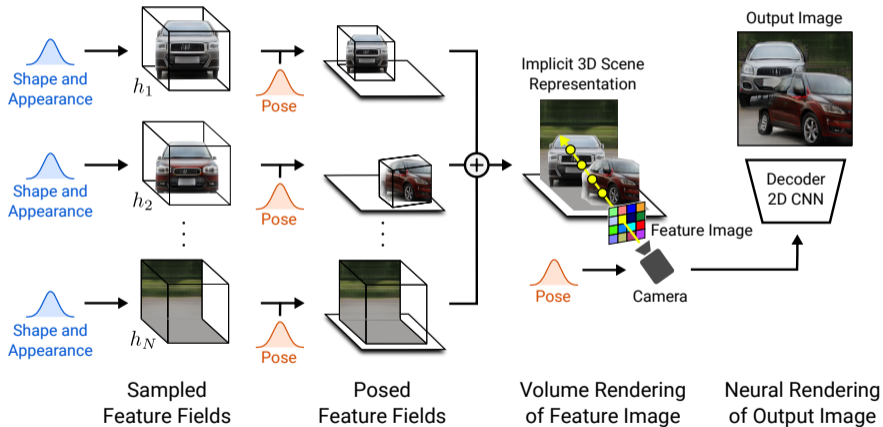
# GIRAFFE

Perform volume rendering and get feature image.



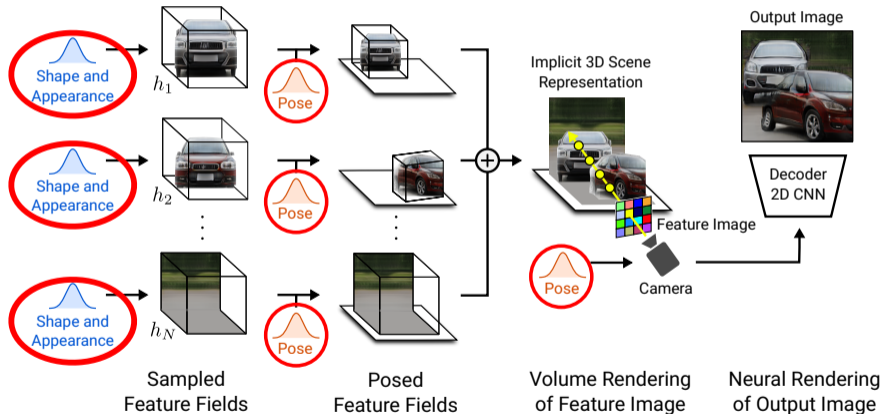
# GIRAFFE

Pass feature image to neural renderer to obtain final output.



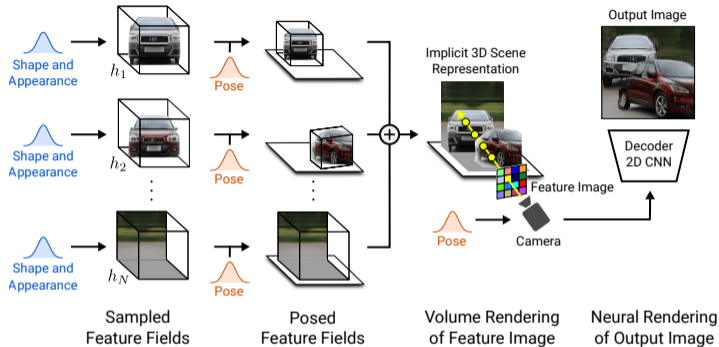
# GIRAFFE

At test time, we can sample individual codes and **control the poses**.





# GIRAFFE



- ▶ We train with adversarial loss **on full image**
- ▶ We volume-render the feature image at  $16 \times 16$  pixels

How well does it work?

# GIRAFFE

We compare object translation for a 2D-based GAN (left) and our method (right):



# GIRAFFE

We can perform more complex operations like circular translations



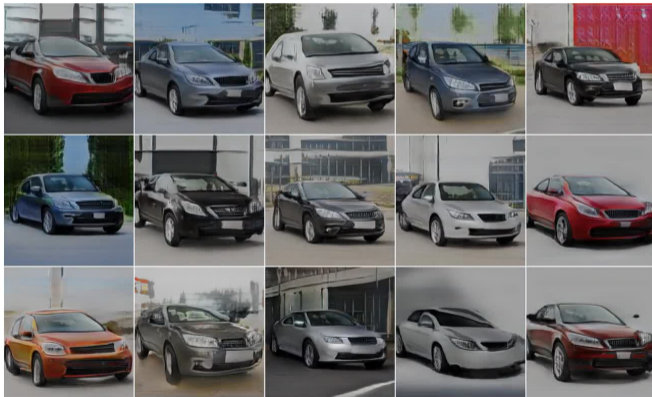
# GIRAFFE

We can add more objects at test time (trained on two-object)



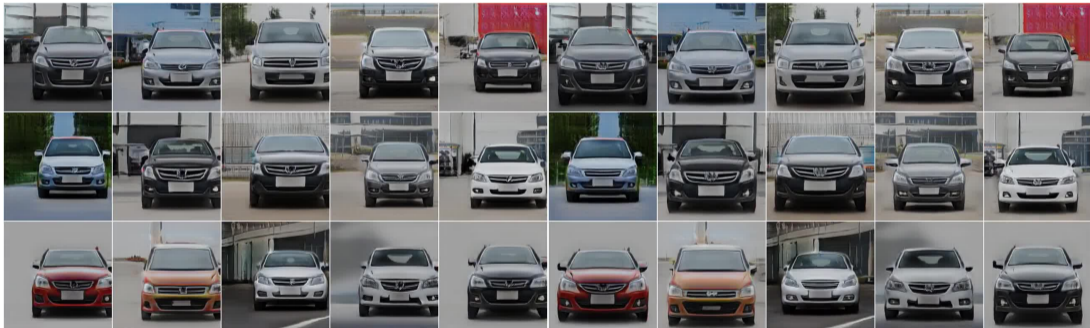
# GIRAFFE

We can rotate the object



# GIRAFFE

We can translate the object



# GIRAFFE

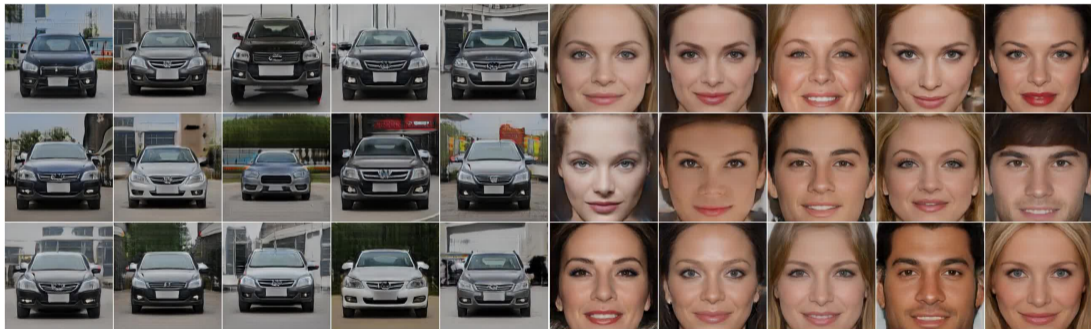
We can change the object shape





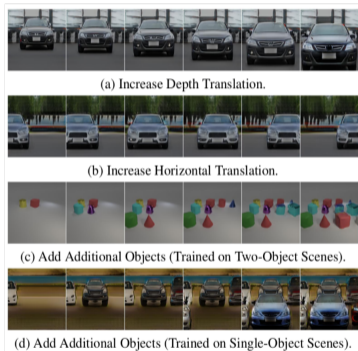
# GIRAFFE

We can change the object appearance



# GIRAFFE

We can generate out-of-distribution samples



How can we scale to more complex camera distributions?

# CAMPARI

GRAF, GIRAFFE:

- ▶ Learn a 3D-aware image generator with uniform prior on camera distributions
- ▶ Requires careful tuning and results degrade if they do not match the data distribution

# CAMPARI

GRAF, GIRAFFE:

- ▶ Learn a 3D-aware image generator with uniform prior on camera distributions
- ▶ Requires careful tuning and results degrade if they do not match the data distribution

CAMPARI:

- ▶ Learn a 3D aware image generator and a **camera generator** jointly.

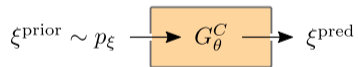
# CAMPARI

Sample prior camera  $\xi^{\text{prior}} \sim p_{\xi}$ .

$$\xi^{\text{prior}} \sim p_{\xi}$$

# CAMPARI

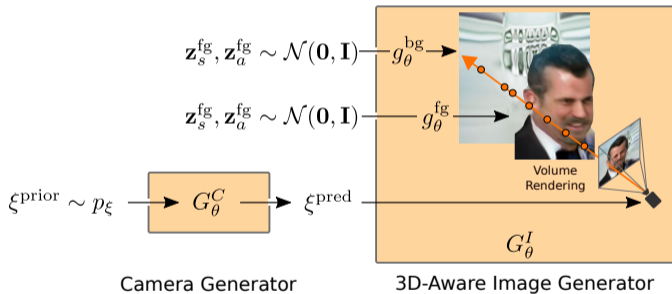
Pass  $\xi^{\text{prior}}$  to camera generator  $G_{\theta}^C$  and obtain predicted camera  $\xi^{\text{pred}}$ .



Camera Generator

# CAMPARI

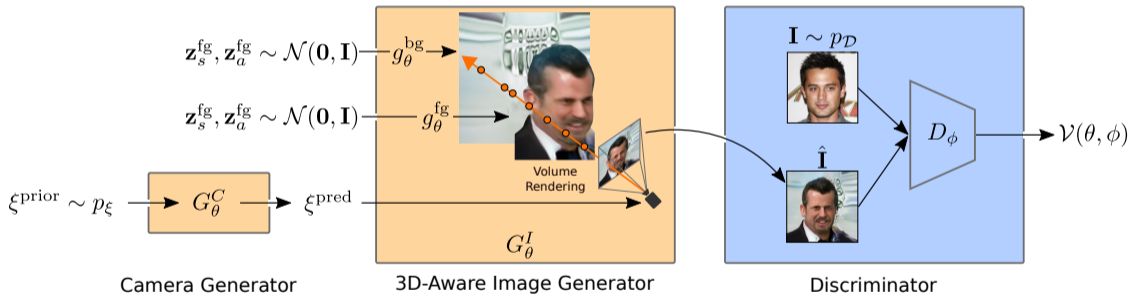
Pass  $\xi^{\text{pred}}$  and sampled FG / BG latent codes to 3D-aware image generator





# CAMPARI

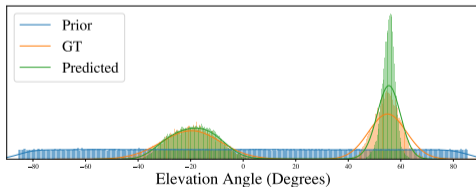
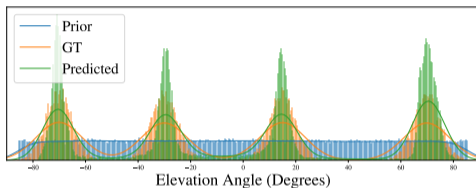
Train entire method with GAN objective (similar to GRAF, GIRAFFE)



How well does it work?

# CAMPARI

CAMPARI learns to match the GT distribution for synthetic datasets



# CAMPARI

## Results on CelebA



(a) Rotation for GRAF [58] (Camera Parameters Tuned)



(b) Rotation for GRAF [58] (Camera Parameters not Tuned)



(c) Rotation for Ours (No Tuning Required)

# Summary

## Summary

- ▶ Occupancy Networks: Represent Surfaces **implicitly** as the decision boundary of a neural network

# Summary

## Summary

- ▶ Occupancy Networks: Represent Surfaces **implicitly** as the decision boundary of a neural network
- ▶ This **coordinate-based neural representation** can also be used for other applications ranging from light fields to motion and more

# Summary

## Summary

- ▶ Occupancy Networks: Represent Surfaces **implicitly** as the decision boundary of a neural network
- ▶ This **coordinate-based neural representation** can also be used for other applications ranging from light fields to motion and more
- ▶ Neural Radiance Fields: Represent scenes as volumes which can be used for view synthesis

# Summary

## Summary

- ▶ Occupancy Networks: Represent Surfaces **implicitly** as the decision boundary of a neural network
- ▶ This **coordinate-based neural representation** can also be used for other applications ranging from light fields to motion and more
- ▶ Neural Radiance Fields: Represent scenes as volumes which can be used for view synthesis
- ▶ Generative Radiance Fields: **3D-aware generative model** of single-object scenes



# Summary

## Summary

- ▶ Occupancy Networks: Represent Surfaces **implicitly** as the decision boundary of a neural network
- ▶ This **coordinate-based neural representation** can also be used for other applications ranging from light fields to motion and more
- ▶ Neural Radiance Fields: Represent scenes as volumes which can be used for view synthesis
- ▶ Generative Radiance Fields: **3D-aware generative model** of single-object scenes
- ▶ GIRAFFE: Generative Model of **complex, multi-object scenes**

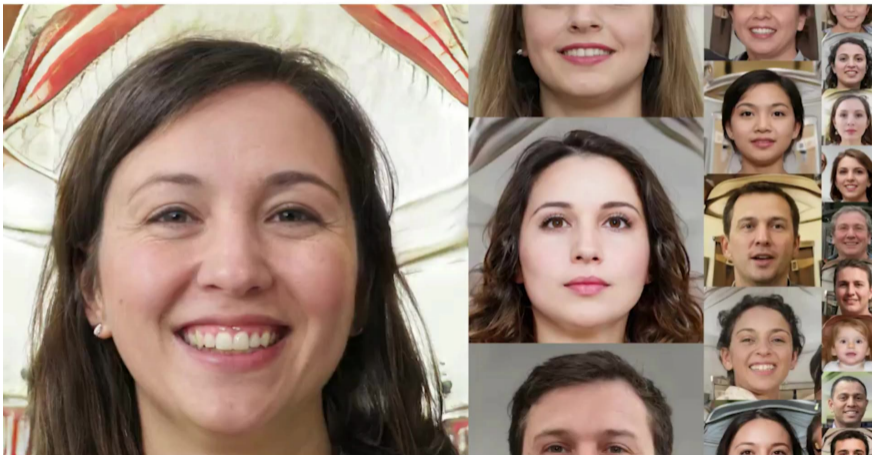
# Summary

## Summary

- ▶ Occupancy Networks: Represent Surfaces **implicitly** as the decision boundary of a neural network
- ▶ This **coordinate-based neural representation** can also be used for other applications ranging from light fields to motion and more
- ▶ Neural Radiance Fields: Represent scenes as volumes which can be used for view synthesis
- ▶ Generative Radiance Fields: **3D-aware generative model** of single-object scenes
- ▶ GIRAFFE: Generative Model of **complex, multi-object scenes**
- ▶ CAMPARI: Generative Model for **more complex camera distributions**

# Summary

This research is very activate and leads to state-of-the-art results:



Thank you!

See <https://m-niemeyer.github.io/> for more information!