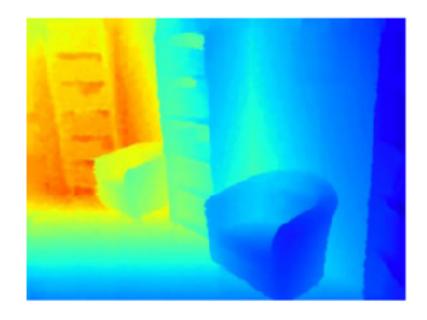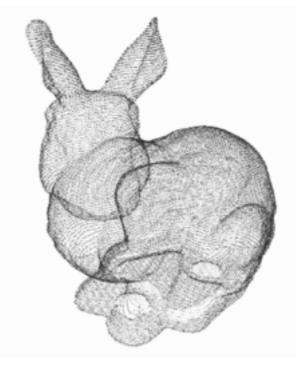# Deep Learning in 3D

## CSE P576
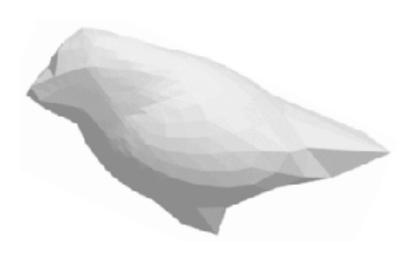
Dr. Matthew Brown

# Deep Learning in 3D

- We'll focus on predicting 3D from one or more image
- Supervision: depth, mesh, silhouettes, view supervision
- Representations: Depth, Points, Meshes, Voxels, SDFs
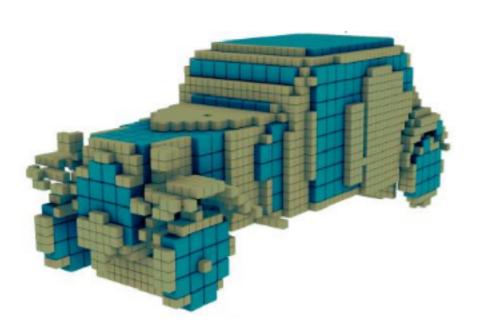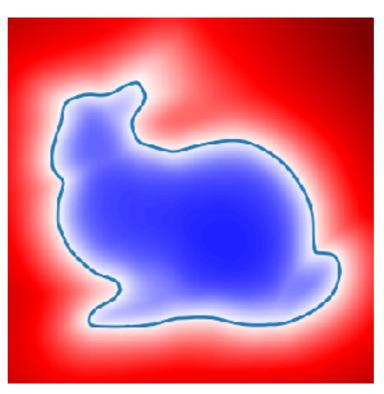- Neural Scene Representation and Rendering
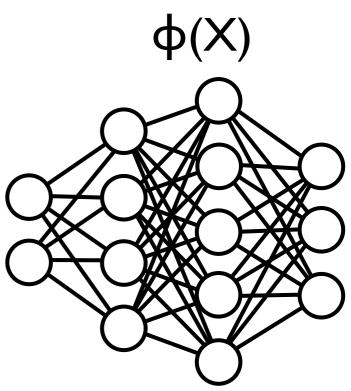
# 3D Representation

- Many ways to represent objects in 3D



$\phi(X)$

# Single-View Depth Estimation

U-Net with skip connections

Direct supervision via Kinect RGB+D

Loss, e.g., L2

# 2-view Stereo

- Form HxWxD=disparity volume and use 3D convolution



| Extract features at each pixel using 2D CNN | Form volume by sliding features from 2nd image at D disparities | Perform 3D convolution on feature volume | Treat output as disparity cost volume and perform soft argmax |

https://www.youtube.com/watch?v=VtAzDS1NLmo  [ Kendall et al. 2017 ]

# Multi-view Stereo



Compare patches in ref image to plane sweep volumes from other images

Perform intra and inter-volume aggregation of features

[ DeepMVS, Huang et al. 2018 ]   6

# DeepMVS: Results



| Image | Ground Truth | Colmap Filtered | Colmap all | DeepMVS [ Huang et al. 2018 ] |

# 3D Shape Representations: Point Cloud

- Represent shape as a set of P points in 3D space
- (+) Can represent fine structures without huge numbers of points
- ( ) Requires new architecture, losses, etc
- (-) Doesn't explicitly represent the surface of the shape: extracting a mesh for rendering or other applications requires post-processing



Fan et al, "A Point Set Generation Network for 3D Object Reconstruction from a Single Image", CVPR 2017

# Processing Pointcloud Inputs: PointNet

Want to process pointclouds as **sets**: order should not matter

Run MLP on each point

Max-Pool

Fully Connected

**Input pointcloud:**
P x 3

**Point features:**
P x D

**Pooled vector:**
D

**Class score:**
C

Qi et al, "PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation", CVPR 2017
Qi et al, "PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space", NeurIPS 2017

# Generating Pointcloud Outputs

Fully connected branch

**Points**: $P_1 \times 3$

**Input Image**: $3 \times H \times W$

**Image Features**: $C \times H' \times W'$

2D CNN

2D CNN

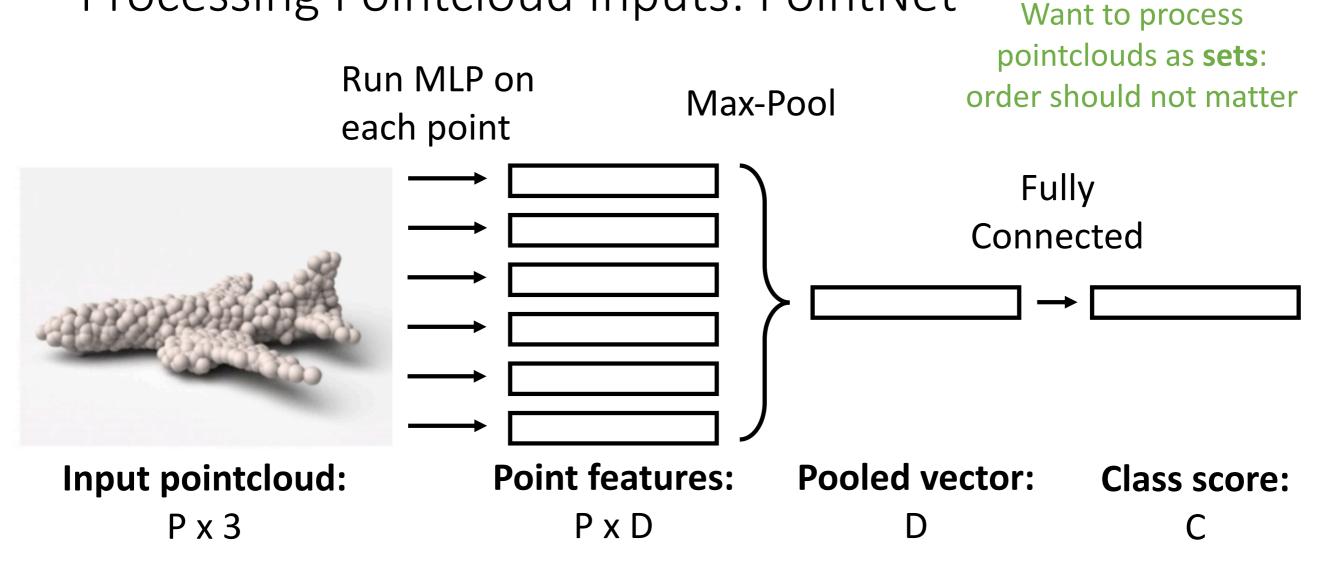**Points**: $(P_2 \times 3) \times H' \times W'$

Convolutional branch

**Pointcloud**: $(P_1 + H'W'P_2) \times 3$

Fan et al, "A Point Set Generation Network for 3D Object Reconstruction from a Single Image", CVPR 2017

10

# Predicting Point Clouds: Loss Function

We need a (differentiable) way to compare pointclouds **as sets**!

**Chamfer distance** is the sum of L2 distance to each point's nearest neighbor in the other set

$$d_{CD}(S_1, S_2) = \sum_{x \in S_1} \min_{y \in S_2} \|x - y\|_2^2 + \sum_{y \in S_2} \min_{x \in S_1} \|x - y\|_2^2$$

Fan et al, "A Point Set Generation Network for 3D Object Reconstruction from a Single Image", CVPR 2017

# Predicting Point Clouds: Loss Function

We need a (differentiable) way to compare pointclouds **as sets**!

**Chamfer distance** is the sum of L2 distance to each point's nearest neighbor in the other set

$$d_{CD}(S_1, S_2) = \sum_{x \in S_1} \min_{y \in S_2} \|x - y\|_2^2 + \sum_{y \in S_2} \min_{x \in S_1} \|x - y\|_2^2$$



Fan et al, "A Point Set Generation Network for 3D Object Reconstruction from a Single Image", CVPR 2017
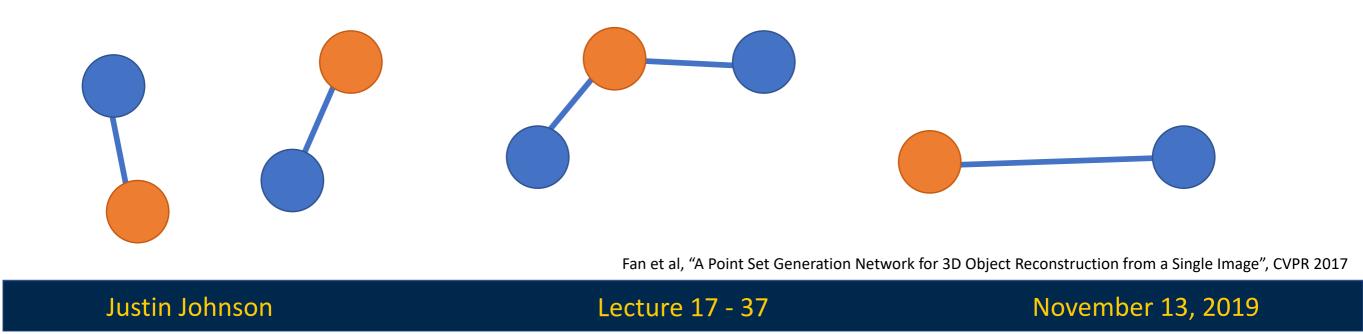
12

# Predicting Point Clouds: Loss Function

We need a (differentiable) way to compare pointclouds **as sets**!

**Chamfer distance** is the sum of L2 distance to each point's nearest neighbor in the other set

$$d_{CD}(S_1, S_2) = \sum_{x \in S_1} \min_{y \in S_2} \|x - y\|_2^2 + \sum_{y \in S_2} \min_{x \in S_1} \|x - y\|_2^2$$
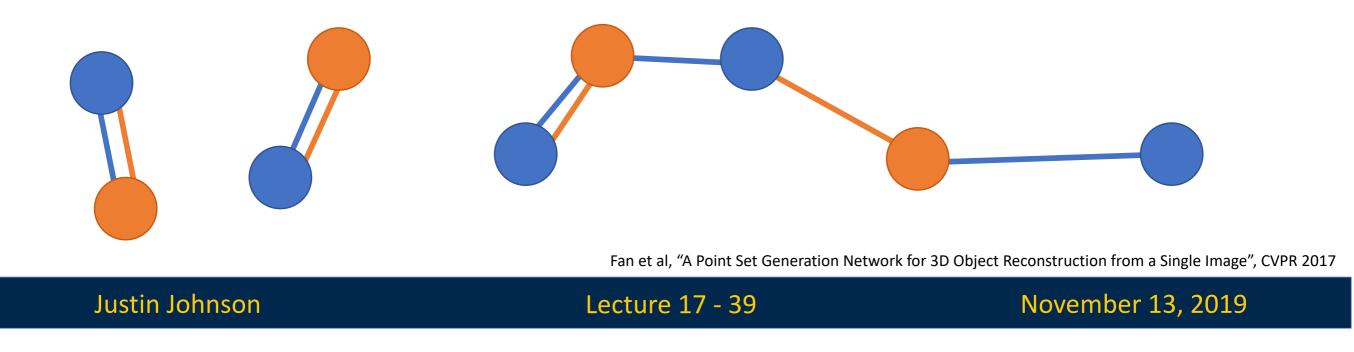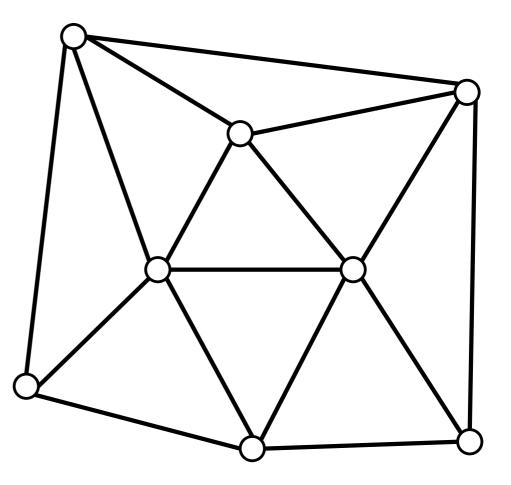
Fan et al, "A Point Set Generation Network for 3D Object Reconstruction from a Single Image", CVPR 2017

ICP-like distance function

# 3D Shape Representations: Triangle Mesh

Represent a 3D shape as a set of triangles

**Vertices**: Set of V points in 3D space

**Faces**: Set of triangles over the vertices

(+) Standard representation for graphics

(+) Explicitly represents 3D shapes

14

# 3D Shape Representations: Triangle Mesh

Represent a 3D shape as a set of triangles

**Vertices**: Set of V points in 3D space

**Faces**: Set of triangles over the vertices

(+) Standard representation for graphics
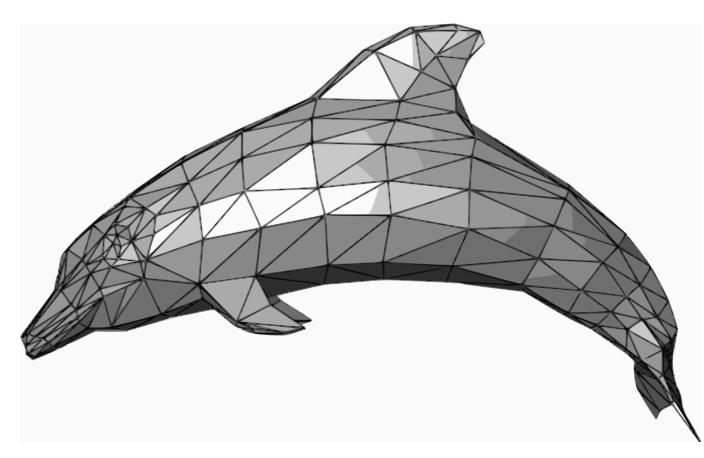
(+) Explicitly represents 3D shapes

(+) Adaptive: Can represent flat surfaces very efficiently, can allocate more faces to areas with fine detail

15

# 3D Shape Representations: Triangle Mesh

Represent a 3D shape as a set of triangles
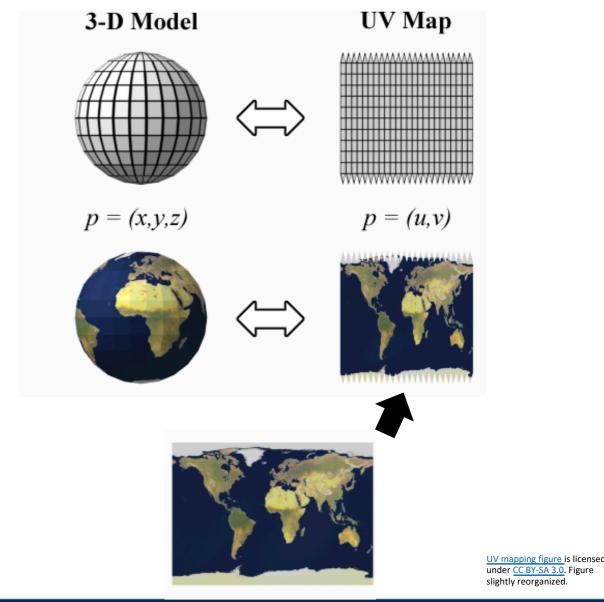
**Vertices**: Set of V points in 3D space

**Faces**: Set of triangles over the vertices

(+) Standard representation for graphics

(+) Explicitly represents 3D shapes

(+) Adaptive: Can represent flat surfaces very efficiently, can allocate more faces to areas with fine detail

(+) Can attach data on verts and interpolate over the whole surface: RGB colors, texture coordinates, normal vectors, etc.

**3-D Model**

**UV Map**

$p = (x,y,z)$

$p = (u,v)$



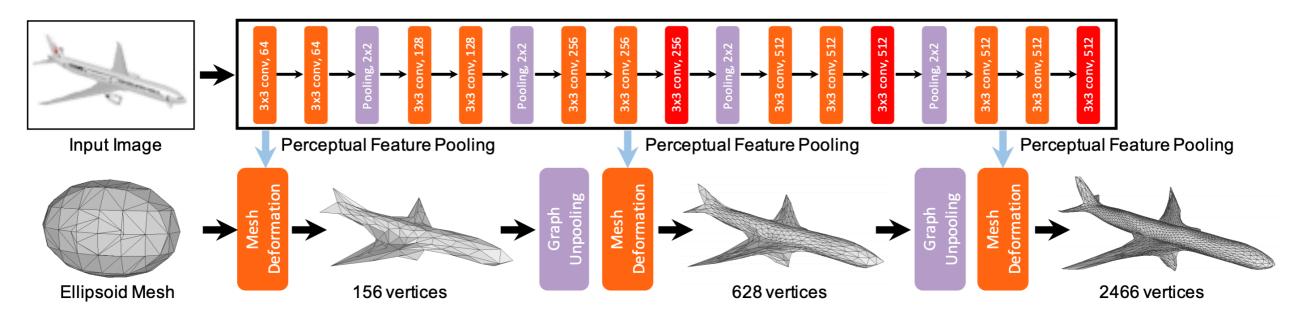UV mapping figure is licensed under CC BY-SA 3.0. Figure slightly reorganized.

16

# Predicting Meshes: Pixel2Mesh

**Input**: Single RGB
Image of an object

**Key ideas**:
Iterative Refinement
Graph Convolution
Vertex Aligned-Features
Chamfer Loss Function

**Output**: Triangle
mesh for the object



Wang et al, "Pixel2Mesh: Generating 3D Mesh Models from Single RGB Images", ECCV 2018

## Supervised with ground truth meshes

17

# Predicting Triangle Meshes: Iterative Refinement

**Idea #1: Iterative mesh refinement**

Start from initial ellipsoid mesh

Network predicts offsets for each vertex

Repeat.

## Fixed mesh structure



Ellipsoid Mesh → Mesh Deformation → 156 vertices → Graph Unpooling → Mesh Deformation → 628 vertices → Graph Unpooling → Mesh Deformation → 2466 vertices

Wang et al, "Pixel2Mesh: Generating 3D Mesh Models from Single RGB Images", ECCV 2018
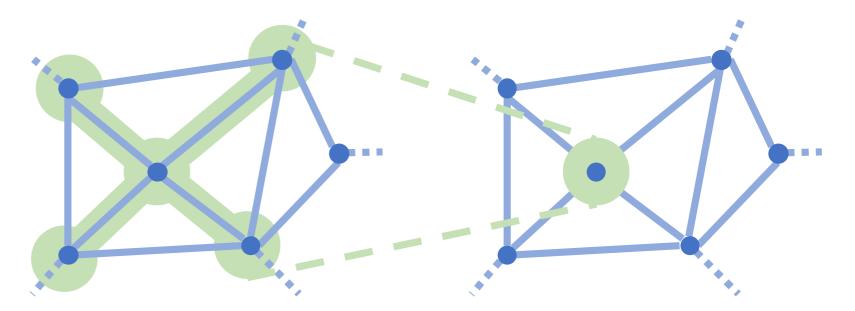
18

# Predicting Triangle Meshes: Graph Convolution

$$f'_i = W_0 f_i + \sum_{j \in \mathcal{N}(i)} W_1 f_j$$

Vertex $v_i$ has feature $f_i$

New feature $f'_i$ for vertex vi depends on feature of neighboring vertices N(i)

Use same weights W0 and W1 to compute all outputs



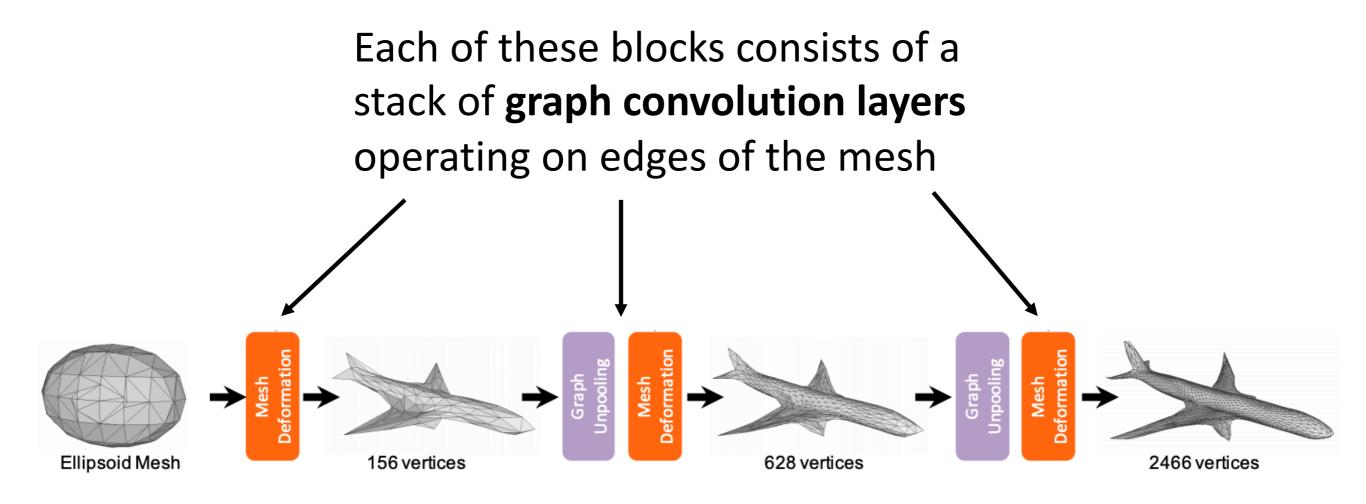**Input**: Graph with a feature vector at each vertex

**Output**: New feature vector for each vertex

19

# Predicting Triangle Meshes: Graph Convolution

Each of these blocks consists of a stack of **graph convolution layers** operating on edges of the mesh
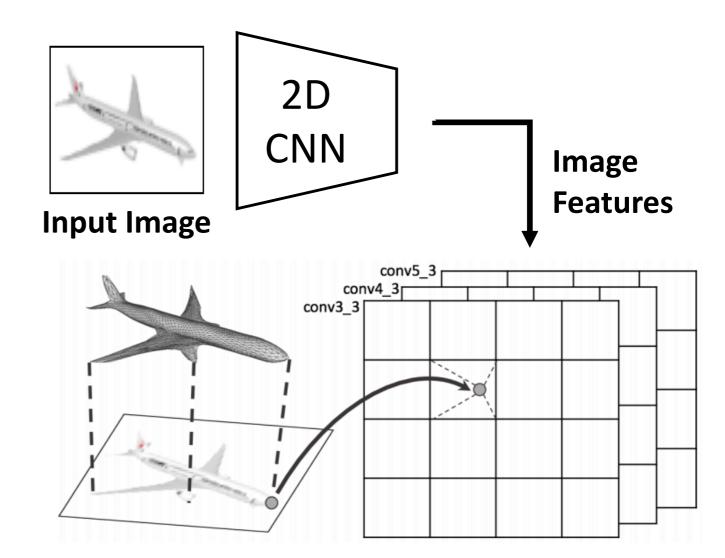


Wang et al, "Pixel2Mesh: Generating 3D Mesh Models from Single RGB Images", ECCV 2018

20

# Predicting Triangle Meshes: Vertex-Aligned Features

**Idea #2**: Aligned vertex features

For each vertex of the mesh:

- Use camera information to project onto image plane

- Use bilinear interpolation to sample a CNN feature

**Input Image**

**2D CNN**

**Image Features**

conv5_3
conv4_3
conv3_3

Wang et al, "Pixel2Mesh: Generating 3D Mesh Models from Single RGB Images", ECCV 2018
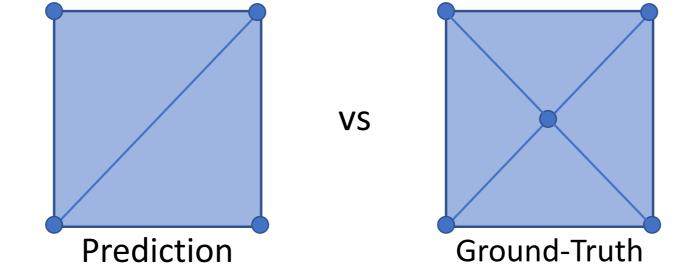
21

# Predicting Meshes: Loss Function

The same shape can be represented with different meshes – how can we define a loss between predicted and ground-truth mesh?

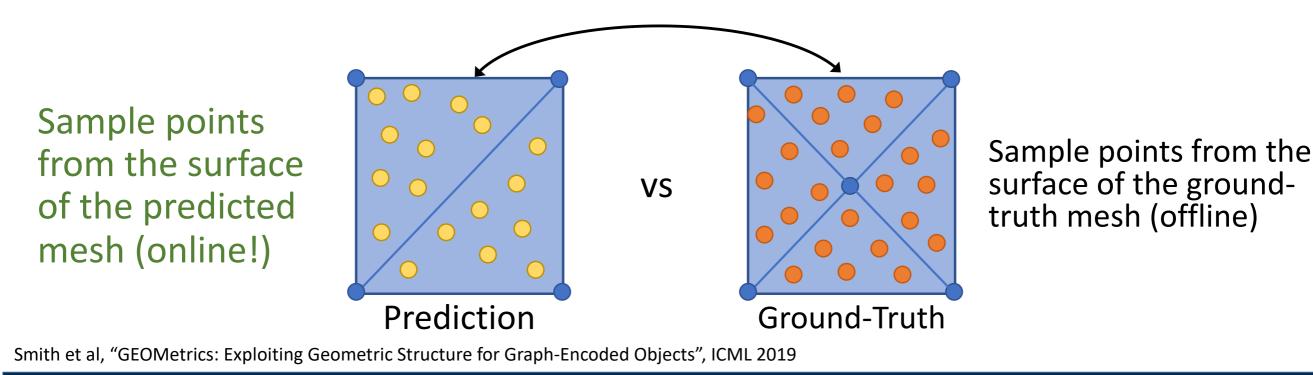**Idea:** Convert meshes to pointclouds, then compute loss



Prediction
vs
Ground-Truth

Wang et al, "Pixel2Mesh: Generating 3D Mesh Models from Single RGB Images", ECCV 2018

# Predicting Meshes: Loss Function

The same shape can be represented with different meshes – how can we define a loss between predicted and ground-truth mesh?

Loss = Chamfer distance between **predicted samples** and **ground-truth samples**

Sample points from the surface of the predicted mesh (online!)

vs

Sample points from the surface of the ground-truth mesh (offline)

Prediction

Ground-Truth

Smith et al, "GEOMetrics: Exploiting Geometric Structure for Graph-Encoded Objects", ICML 2019

# Predicting Meshes: Pixel2Mesh

**Key ideas**:
Iterative Refinement
Graph Convolution
Vertex Aligned-Features
Chamfer Loss Function

**Input**: Single RGB
Image of an object

**Output**: Triangle
mesh for the object



Input Image

| 3x3 conv, 64 | 3x3 conv, 64 | Pooling, 2x2 | 3x3 conv, 128 | 3x3 conv, 128 | Pooling, 2x2 | 3x3 conv, 256 | 3x3 conv, 256 | 3x3 conv, 256 | Pooling, 2x2 | 3x3 conv, 512 | 3x3 conv, 512 | 3x3 conv, 512 | Pooling, 2x2 | 3x3 conv, 512 | 3x3 conv, 512 | 3x3 conv, 512 |

Perceptual Feature Pooling     Perceptual Feature Pooling     Perceptual Feature Pooling

Ellipsoid Mesh → Mesh Deformation → 156 vertices → Graph Unpooling → Mesh Deformation → 628 vertices → Graph Unpooling → Mesh Deformation → 2466 vertices

Wang et al, "Pixel2Mesh: Generating 3D Mesh Models from Single RGB Images", ECCV 2018

Supervised with ground truth meshes

24

# Category Specific Mesh Reconstruction
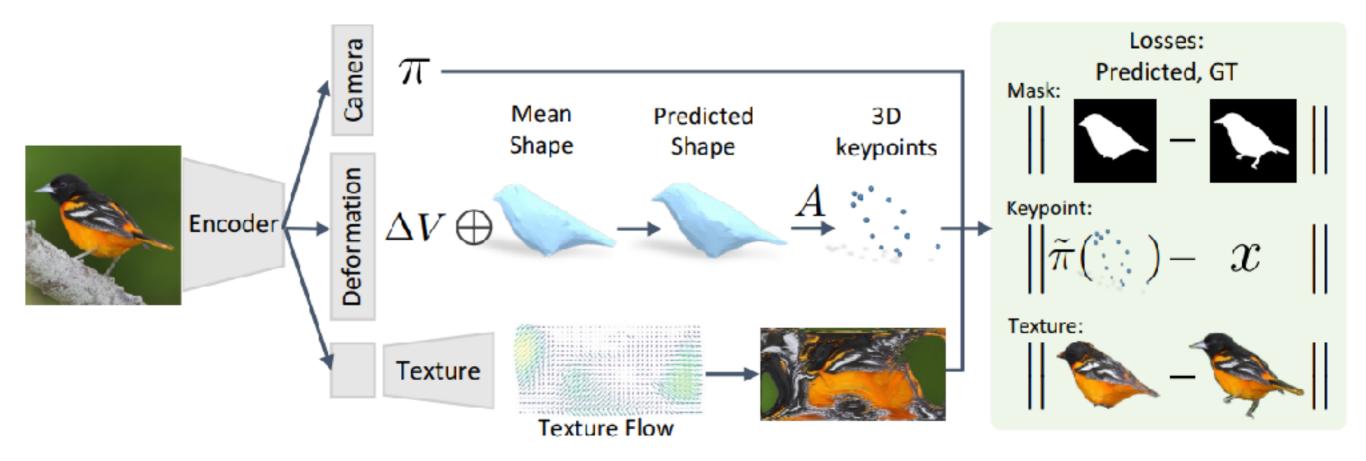
- Can we learn without ground truth meshes?



Given an image, infer mesh, camera, texture

[ Kanazawa et al. 2018 ]  25

Data = Caltech-UCSD birds CUB-200-2011, 6000 images of 200 bird species, + segmentation, 14 semantic keypoints, remove 300 images where num visible keypoints <= 6

# Category Specific Mesh Reconstruction

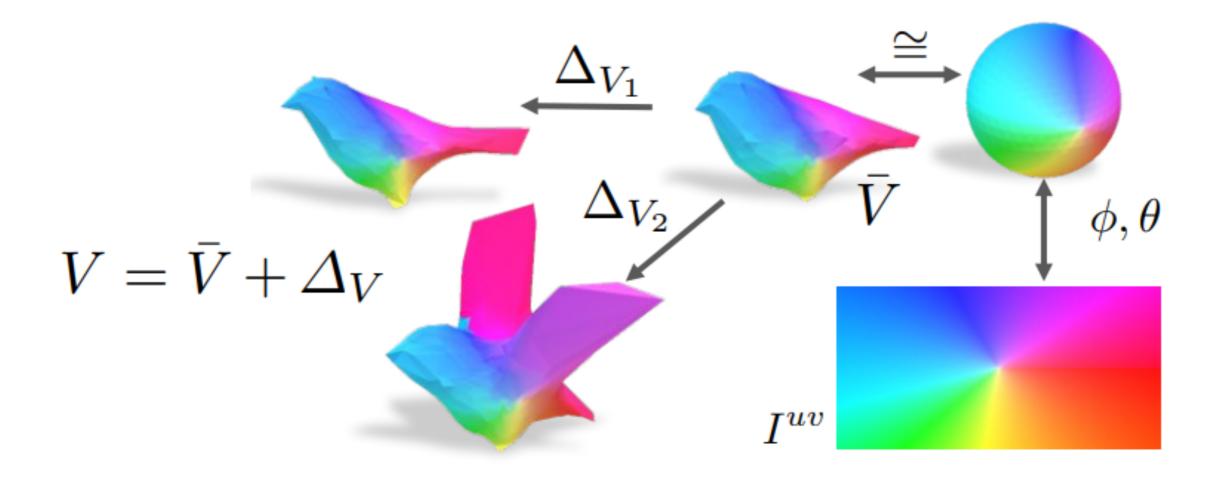- Train a model to predict object mesh (deformation of mean category shape) + camera pose



- Use semantic keypoints and object masks to learn shape (texture not used to learn shape in this implementation)

[ Kanazawa et al. 2018 ] 27

# Mesh Parametrization

- Fixed spherical mesh (subdivided icosahedron) 642 vertices V 1280 faces
- Instances are deformations ΔV of a mean class shape V
- Texture is modelled as RGB colour in spherical coordinates



$$V = \bar{V} + \Delta_V$$

# Keypoints and Projection

- Semantic keypoint positions are modelled as weighted vertex positions

- Matrix A is learned per-class, can be viewed as per-vertex probabilities, with keypoints as the expected value

- Projection $\pi$ is modelled by a camera with translation t, rotation q (quaternion) and scale s



$$v_k = \sum_v A_{k,v} v$$

$$A \in \mathcal{R}_+^{|K| \times |V|}$$

$A \cdot V$ = set of keypoint positions

# Keypoint Projection Loss

- Ensure that keypoints (parametrized as weighted vertex positions) map to the known positions xi. Note: weightings A are per class, vertices V per instance



$$L_{\text{reproj}} = \sum_i ||x_i - \tilde{\pi}_i(AV_i)||_2$$

# Mask Projection Loss

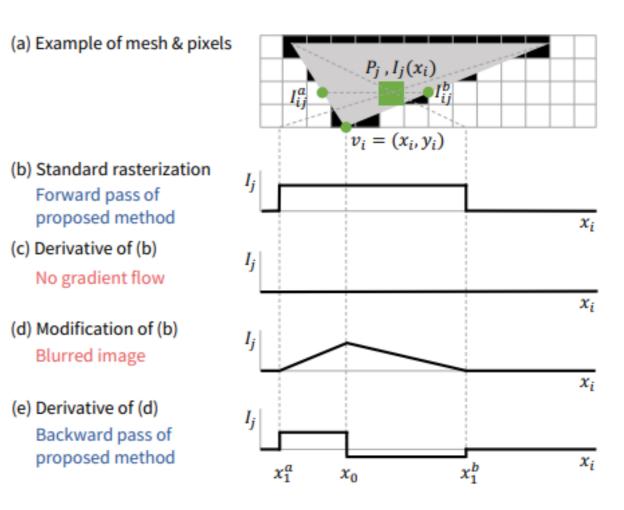- Ensure that the mesh maps to the known silhouette. Note: gradient depends on rendering the mesh



$$L_{\mathrm{mask}} = \sum_i \left\| S_i - \mathcal{R}(V_i, F, \tilde{\pi}_i) \right\|_2$$

S = silhouette, R(.) mesh rendering

# Gradient of Mesh Render

Extend gradient for each pixel inside/
outside triangles with linear ramp



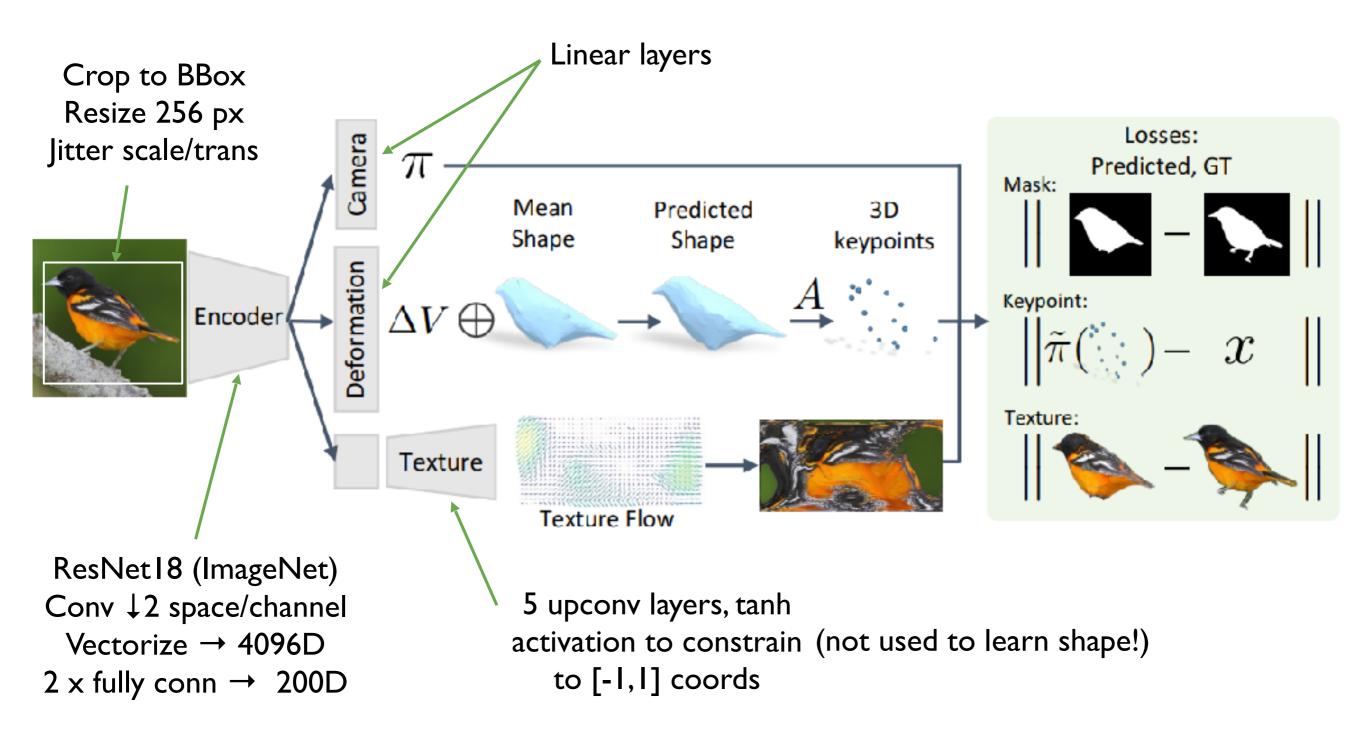[ Neural 3D Mesh Renderer, Kato et al. CVPR 2018 ]

# Texture Representation

- Texture is parametrized as coordinates (flow) of the input image I(u,v)
- → each point on the reference sphere is given a coordinate in the input image
- Latent representation is upconvolved to generate flow I(u,v)
- Loss is Zhang et al. perceptual loss [1] of projected texture
- Note: texture loss is not used to learn shape!



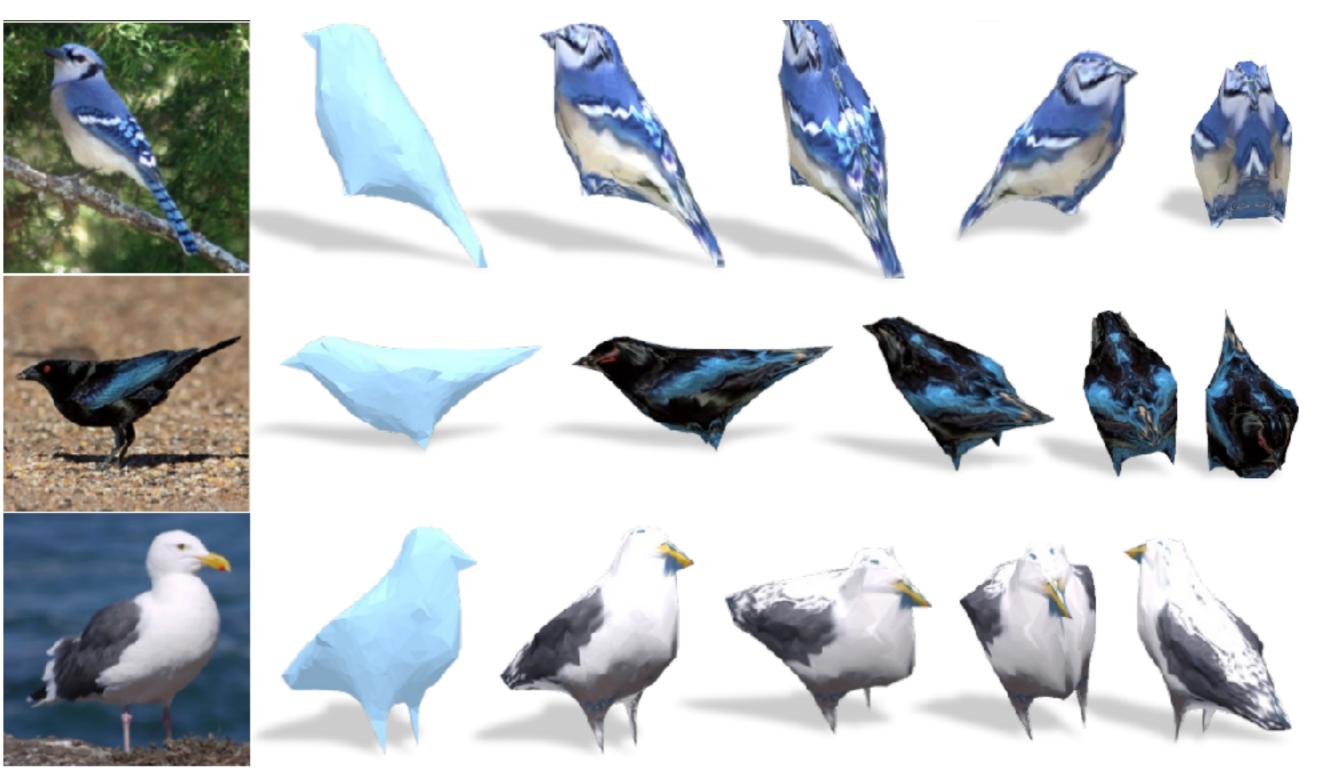[1] The unreasonable effectiveness of deep networks as a perceptual metric. R. Zhang et al. CVPR 2018

Crop to BBox
Resize 256 px
Jitter scale/trans

Linear layers

Camera

$\pi$

Deformation

$\Delta V \oplus$

Mean Shape

Predicted Shape

3D keypoints

$A$

Texture

Texture Flow

Losses:
Predicted, GT

Mask:

$$\left\|\quad - \quad\right\|$$

Keypoint:

$$\left\|\tilde{\pi}(\quad) - x\right\|$$

Texture:

$$\left\|\quad - \quad\right\|$$

ResNet18 (ImageNet)
Conv ↓2 space/channel
Vectorize → 4096D
2 x fully conn → 200D

5 upconv layers, tanh
activation to constrain (not used to learn shape!)
to [-1,1] coords

# SFM Initialization

- In principle, camera $\pi$, mean shape V, instance shape $\Delta$V, keypoint weightings A could be learned from supervised keypoint and silhouette losses
- In practice, the authors initialize cameras $\pi$ and mean shape V via SFM
- Note this involves bundle adjustment / optimization over different birds, so results in fitting an "average" bird model
- The mesh is initialized as the convex hull of keypoint positions, and camera solutions $\pi$_hat are recorded

Initial Mean Shape

# Results

# Results

# Deformation Modes

- Mean and first 3 PCA components of bird shapes



Initial Mean Shape

Learned Mean Shape

# 3D Shape Prediction: Mesh R-CNN

Mask R-CNN:
2D Image -> 2D shapes

**Mesh** R-CNN:
2D Image -> **Triangle Meshes**



He, Gkioxari, Dollár, and Girshick, "Mask R-CNN", ICCV 2017

Gkioxari, Malik, and Johnson, "Mesh R-CNN", ICCV 2019

Detect objects and extract silhouettes

Estimate 3D mesh

39

# 3D Datasets: Object-Centric

uses 3D mesh models from IKEA

## ShapeNet



club chair    cantilever chair    armchair

~50 categories, ~50k 3D CAD models

Standard split has 13 categories, ~44k models, 25 rendered images per model

Many papers show results here

(-) Synthetic, isolated objects; no context

(-) Lots of chairs, cars, airplanes

Chang et al, "ShapeNet: An Information-Rich 3D Model Repository", arXiv 2015
Choy et al, "3D-R2N2: A Unified Approach for Single and Multi-view 3D Object Reconstruction", ECCV 2016

## Pix3D



9 categories, 219 3D models of IKEA furniture aligned to ~17k real images

Some papers train on ShapeNet and show qualitative results here, but use ground-truth segmentation masks

(+) Real images! Context!

(-) Small, partial annotations – only 1 obj/image

Sun et al, "Pix3D: Dataset and Methods for Single-Image 3D Shape Modeling", CVPR 2018

40

# Mesh R-CNN: Task

**Input**: Single RGB image

**Output**:
A set of detected objects
For each object:
- Bounding box
- Category label
- Instance segmentation
- 3D triangle mesh

Mask R-CNN

Mesh head

41

# Mesh R-CNN: Hybrid 3D shape representation

**Mesh deformation** gives good results, but the topology (verts, faces, genus, connected components) fixed by the initial mesh



Ellipsoid Mesh → Mesh Deformation → 156 vertices

**Our approach**: Use voxel predictions to create initial mesh prediction!

# Mesh R-CNN Pipeline

Input image



2D object recognition



3D object meshes



3D object voxels

43

# Mesh R-CNN: ShapeNet Results

# Mesh R-CNN: Pix3D Results

Amodal completion: predict
occluded parts of objects



Box & Mask Predictions

Mesh Predictions

# 3D Shape Representations: Voxels

- Represent a shape with a V x V x V grid of occupancies
- Just like segmentation masks in Mask R-CNN, but in 3D!
- (+) Conceptually simple: just a 3D grid!
- (-) Need high spatial resolution to capture fine structures
- (-) Scaling to high resolutions is nontrivial!



Choy et al, "3D-R2N2: A Unified Approach for Single and Multi-view 3D Object Reconstruction", ECCV 2016

# Processing Voxel Inputs: 3D Convolution



Input:
1 x 30 x 30 x 30

6x6x6 conv
48x13x13x13

5x5x5 conv
160x5x5x5

4x4x4 conv
512x2x2x2

FC
Layer

Class
Scores

Train with classification loss

Wu et al, "3D ShapeNets: A Deep Representation for Volumetric Shapes", CVPR 2015

(for classification of a voxel grid)

# Generating Voxel Shapes: 3D Convolution



**3D CNN**

Input image:
3 x 112 x 112

2D Features:
C x H x W

3D Features:
C' x D' x H' x W'

Voxels:
1 x V x V x V

Train with per-voxel cross-entropy loss

Choy et al, "3D-R2N2: A Unified Approach for Single and Multi-view 3D Object Reconstruction", ECCV 2016

# Voxel Problems: Memory Usage

Storing $1024^3$ voxel grid takes 4GB of memory!

Voxel memory usage (V x V x V float32 numbers)

# Scaling Voxels: Oct-Trees

## Use voxel grids with heterogenous resolution!



Tatarchenko et al, "Octree Generating Networks: Efficient Convolutional Architectures for High-resolution 3D Outputs", ICCV 2017
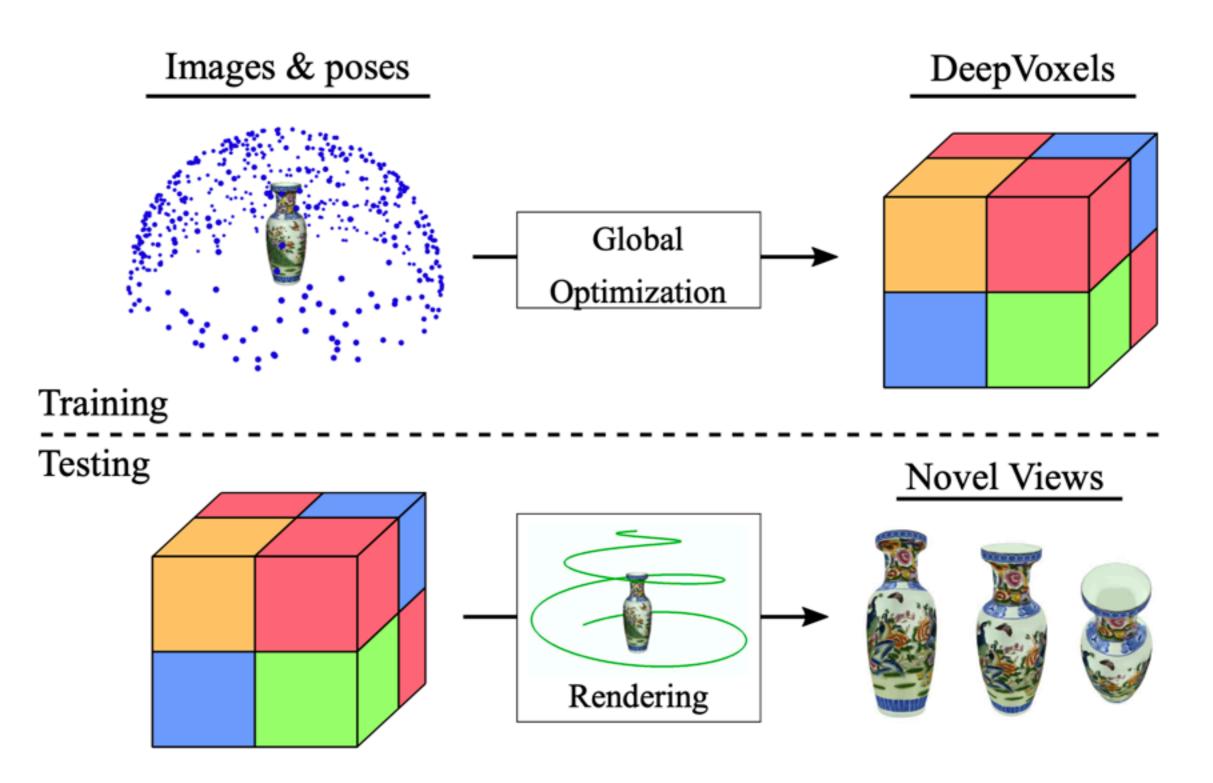
# Self-supervised Scene Representation Learning



Observations → Model, e.g., depth maps, mesh, voxels, SDF → Differentiable Render → Re-Rendered Observations

Image Loss

# Self-supervised Scene Representation Learning

# DeepVoxels

## Embedding vector per voxel



Observations

Neural Scene
Representation

Neural Renderer

Re-Rendered
Observations

Image Loss

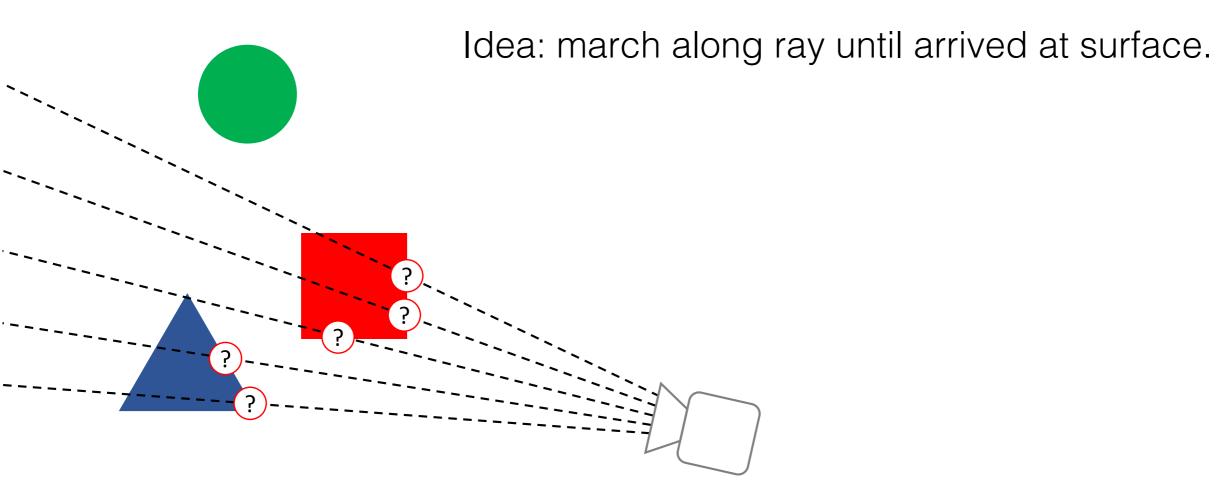Scene represented as an embedding vector per 3D point

# DeepVoxels



[ Sitzmann et al. 2019 ]  54

Ground Truth          Worral et al          pix2pix          DeepVoxels

# Scene Representation Networks

Fully connected network φ=f(X)



Observations

Neural Scene Representation

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

Neural Renderer

Re-Rendered Observations

Image Loss

Scene represented as an embedding vector per 3D point

# Image Regression

- Networks that operate on coordinates to generate image representations are sometimes called "Compositional Pattern Producing Networks" (CPPNs)

$$(R, G, B) = \phi(x, y)$$



[ A. Karpathy ConvNetJS Image Regression demo ]

# Scene Representation Networks

Neural Scene
Representation

Observations

Re-Rendered
Observations

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

Neural Renderer

Image Loss

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.

Iteration 0

# Neural Renderer Step 1: Intersection Testing.

Iteration 1

Neural Renderer Step 1: Intersection Testing.
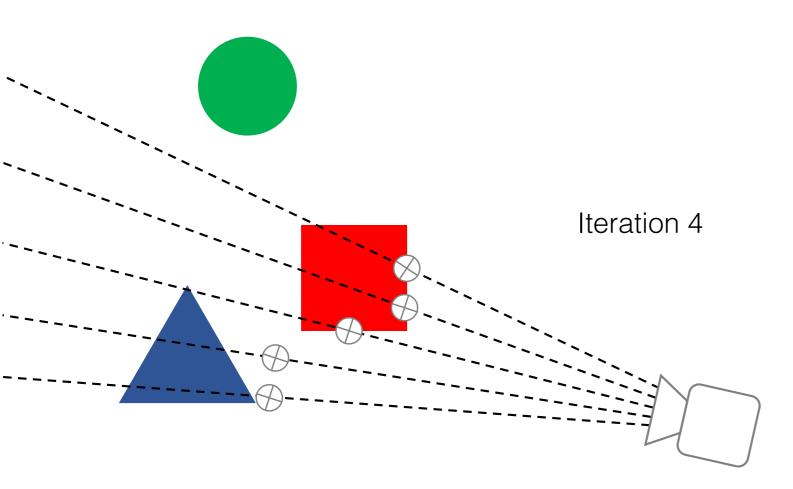
Iteration 2

# 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

Per-pixel Embedding →RGB

# 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

# View Synthesis: Shapenet Cars

- Train using 50 observations per object, known cameras

Observation: <u>Single</u> Image



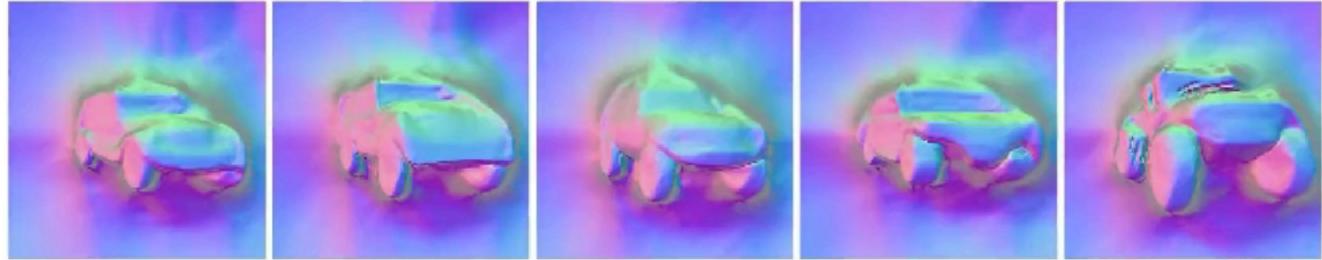Model Output: Novel Views



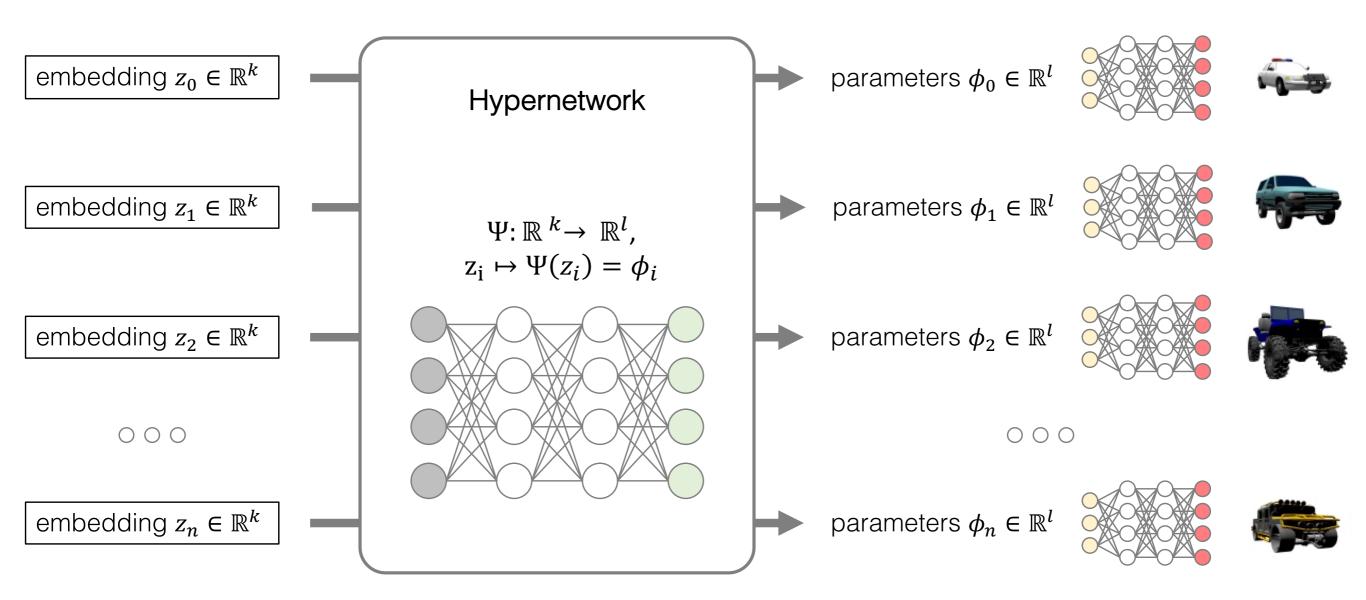Model Output: Geometry (unsupervised)

# Sampling at arbitrary resolutions



512x512

32x32

64x64

128x128

256x256

Surface Normals

RGB

Can render scene at any resolution ɸ=f(X)

# Each scene represented by its own SRN.



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

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

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

$\circ \circ \circ$

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

Hypernetwork

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

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

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

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

$\circ \circ \circ$

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

# Latent Code Interpolation

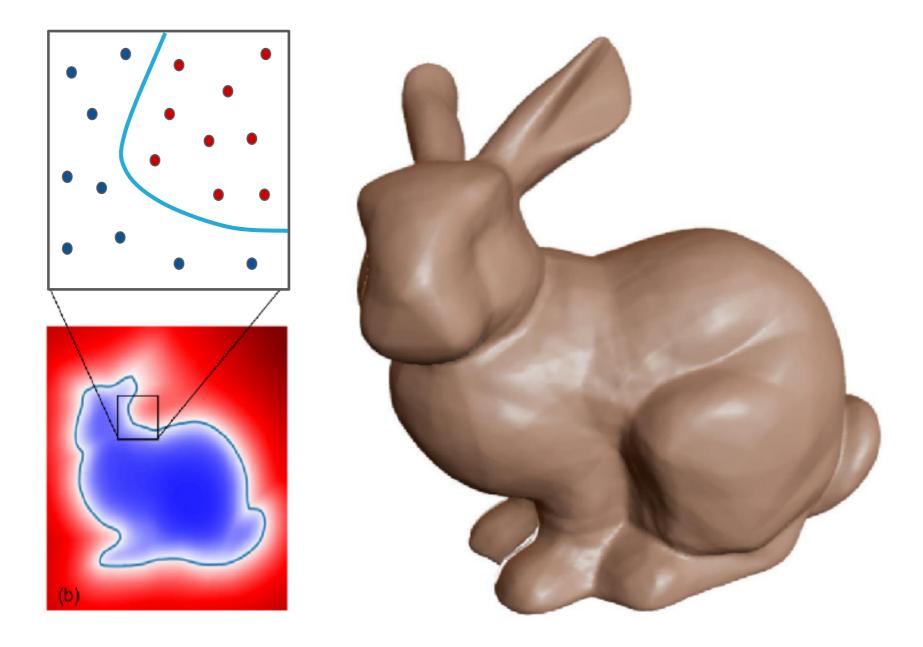- Interpolated latent codes give meaningful scenes

# DeepSDF: Learning Continuous SDFs for Shape Representation

Jeong Joon Park, Peter Florence, Julian Straub,
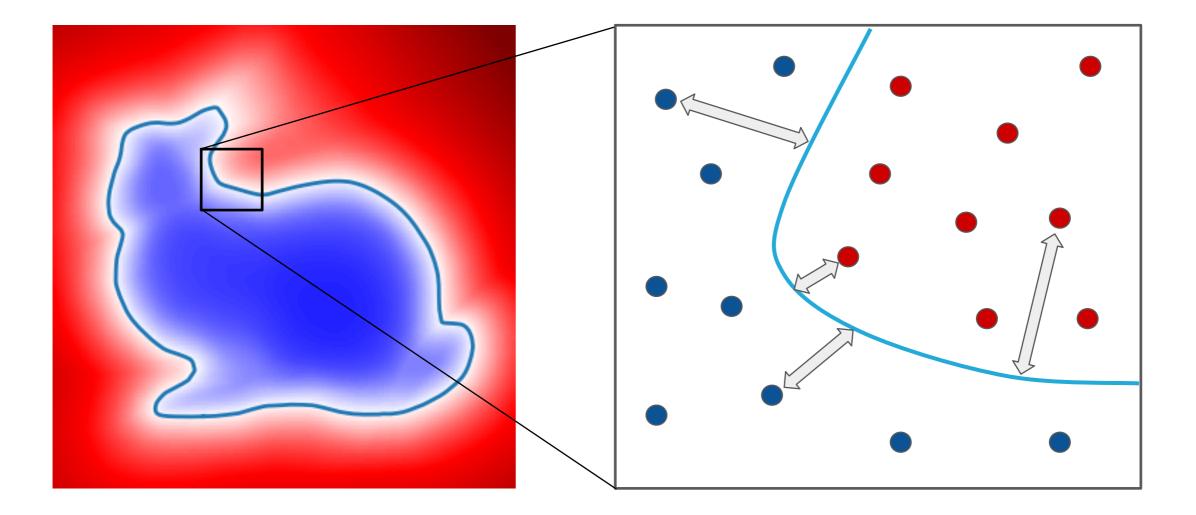Richard Newcombe, Steven Lovegrove
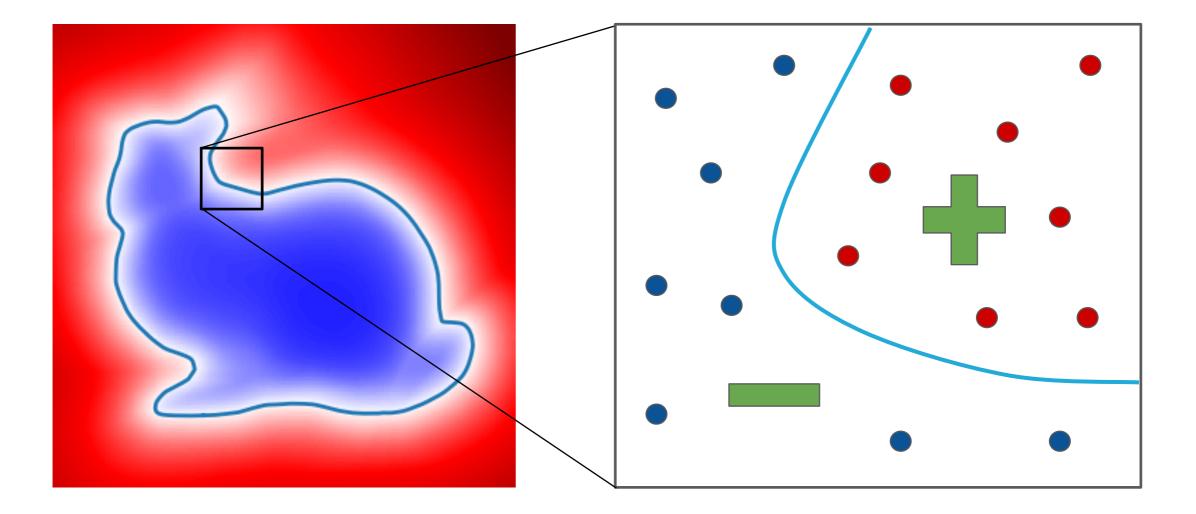
CVPR 2019

# DeepSDF

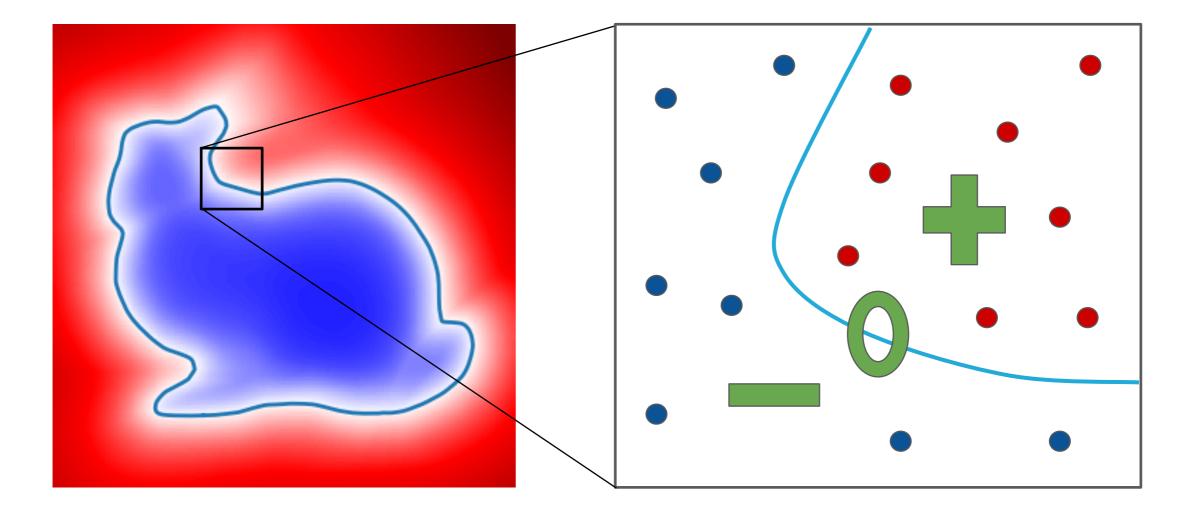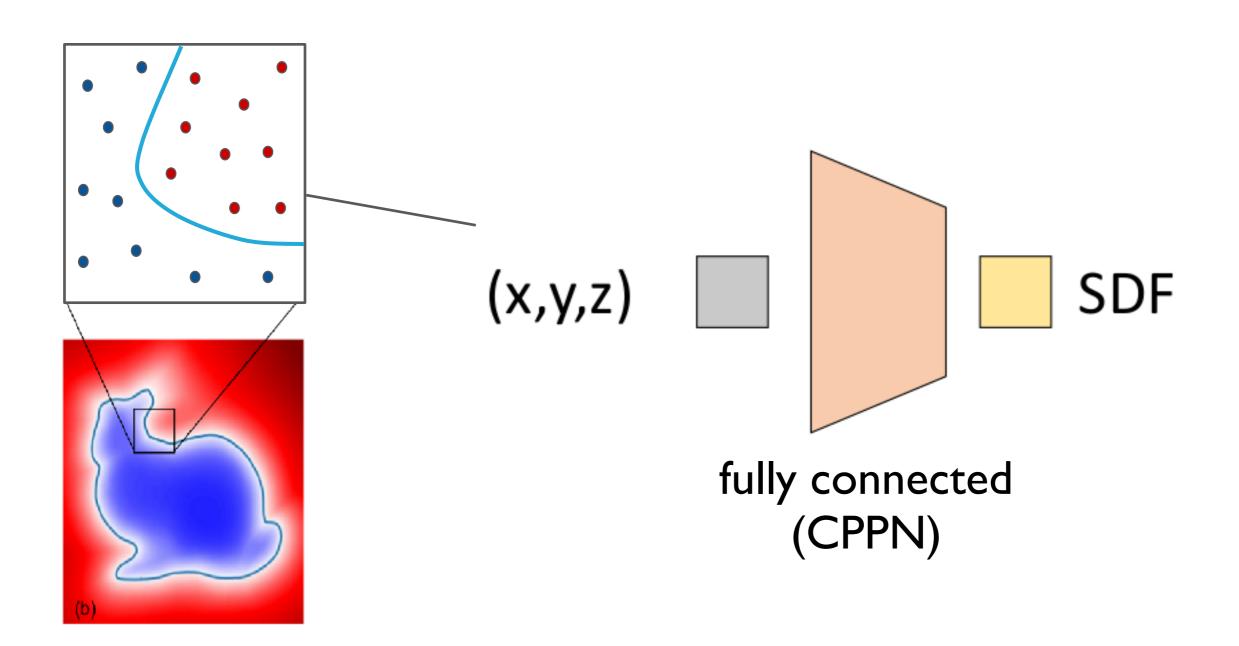- CPPN for signed distance function, SDF=f(X)



[ Slides: Jeong Joon Park ]

# Signed Distance Function

# Signed Distance Function

# Signed Distance Function



[ Slides: Jeong Joon Park ]

# SDF Regression



$(x,y,z)$

SDF

fully connected
(CPPN)
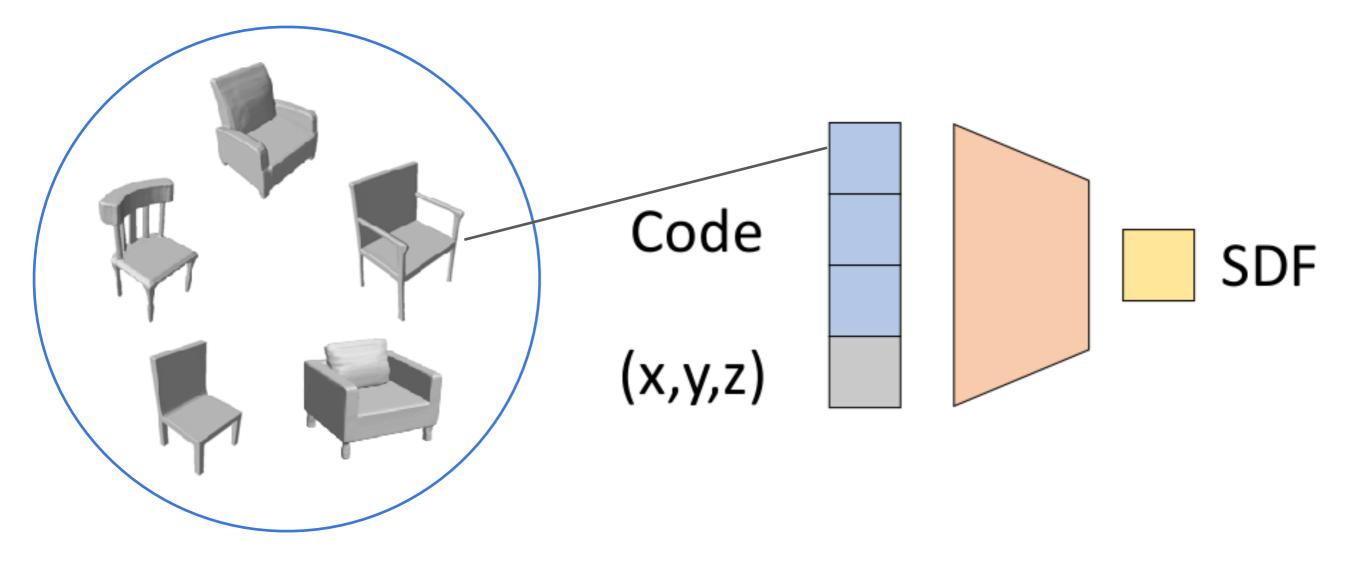
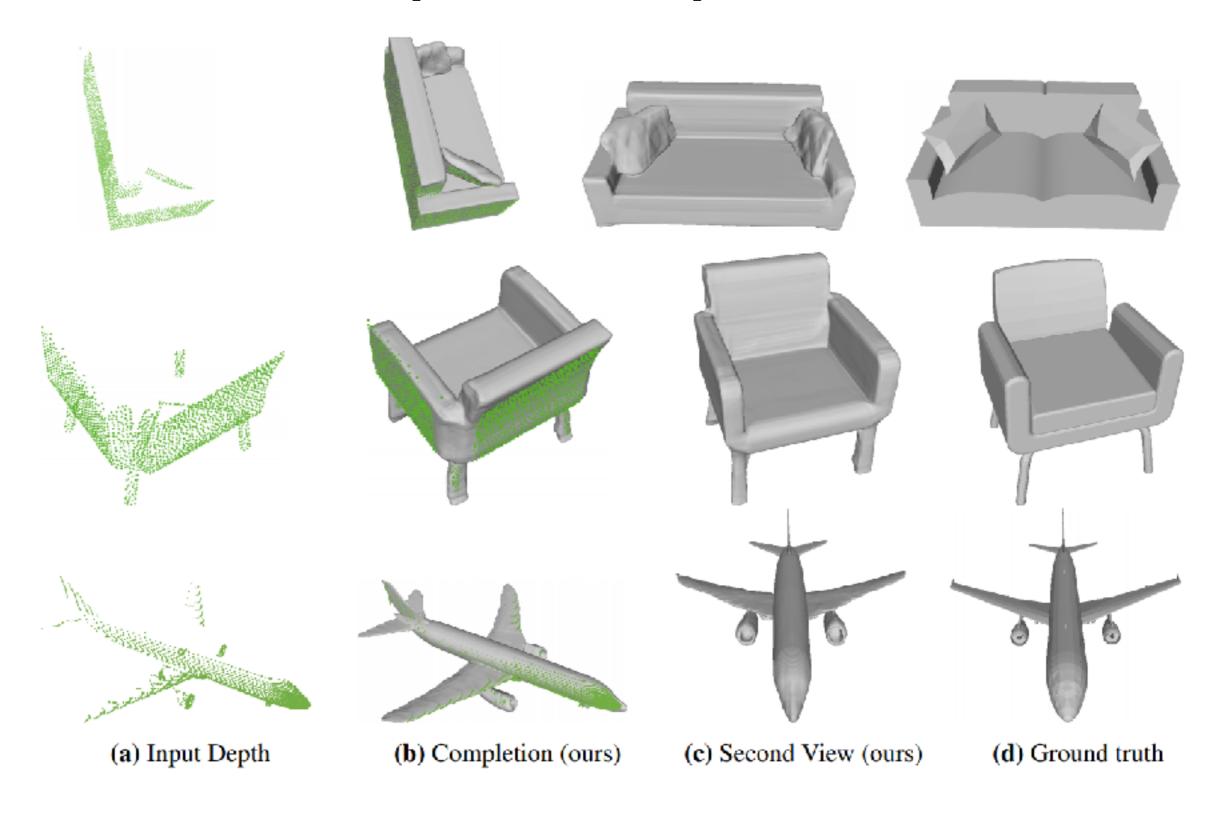Estimate parameters of fully connected net f(X) to fit known SDF

# Shape Modelling

## Coding Multiple Shapes



Assign random codes to each training object,
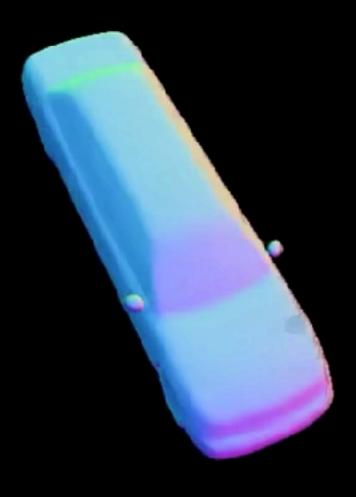optimise network parameters to fit known 3D

# Shape Completion



(a) Input Depth    (b) Completion (ours)    (c) Second View (ours)    (d) Ground truth
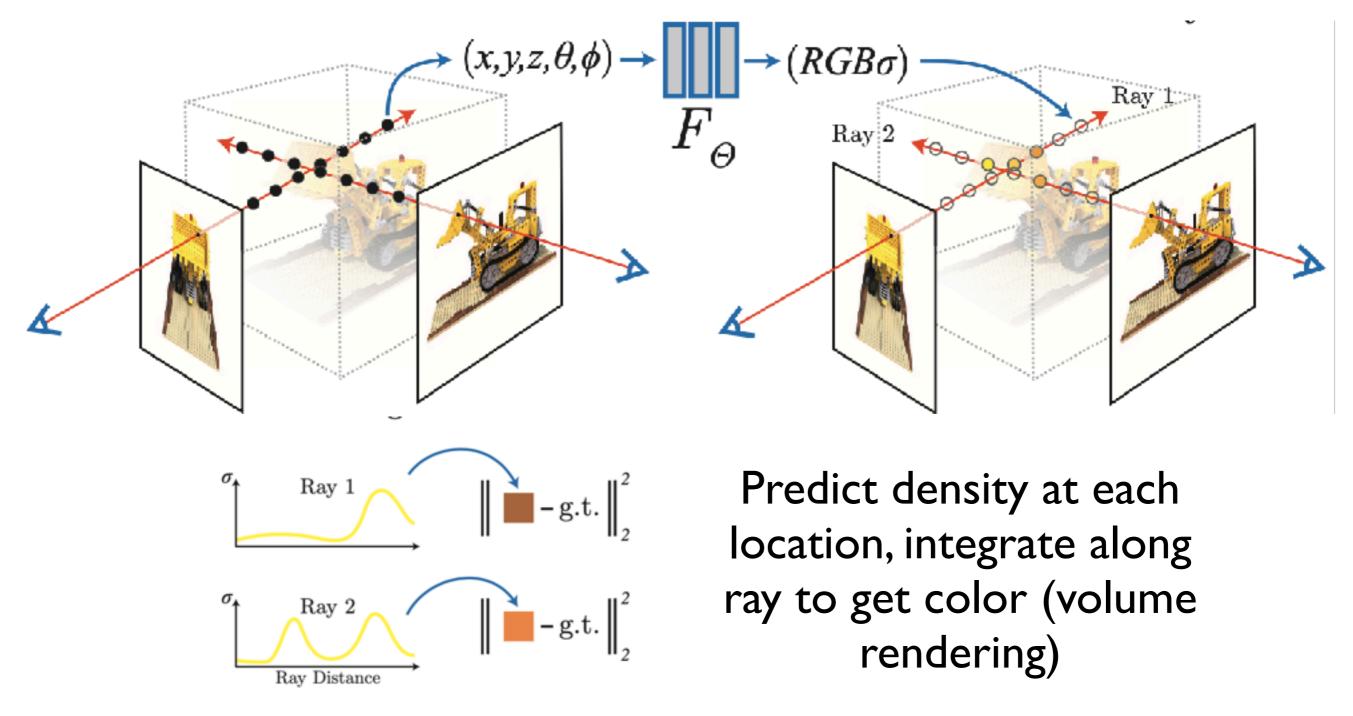
Optimise latent code given partial SDF by backprop to input

Learned Chair Shape Space

Learned Car Shape Space

# Neural Radiance Fields

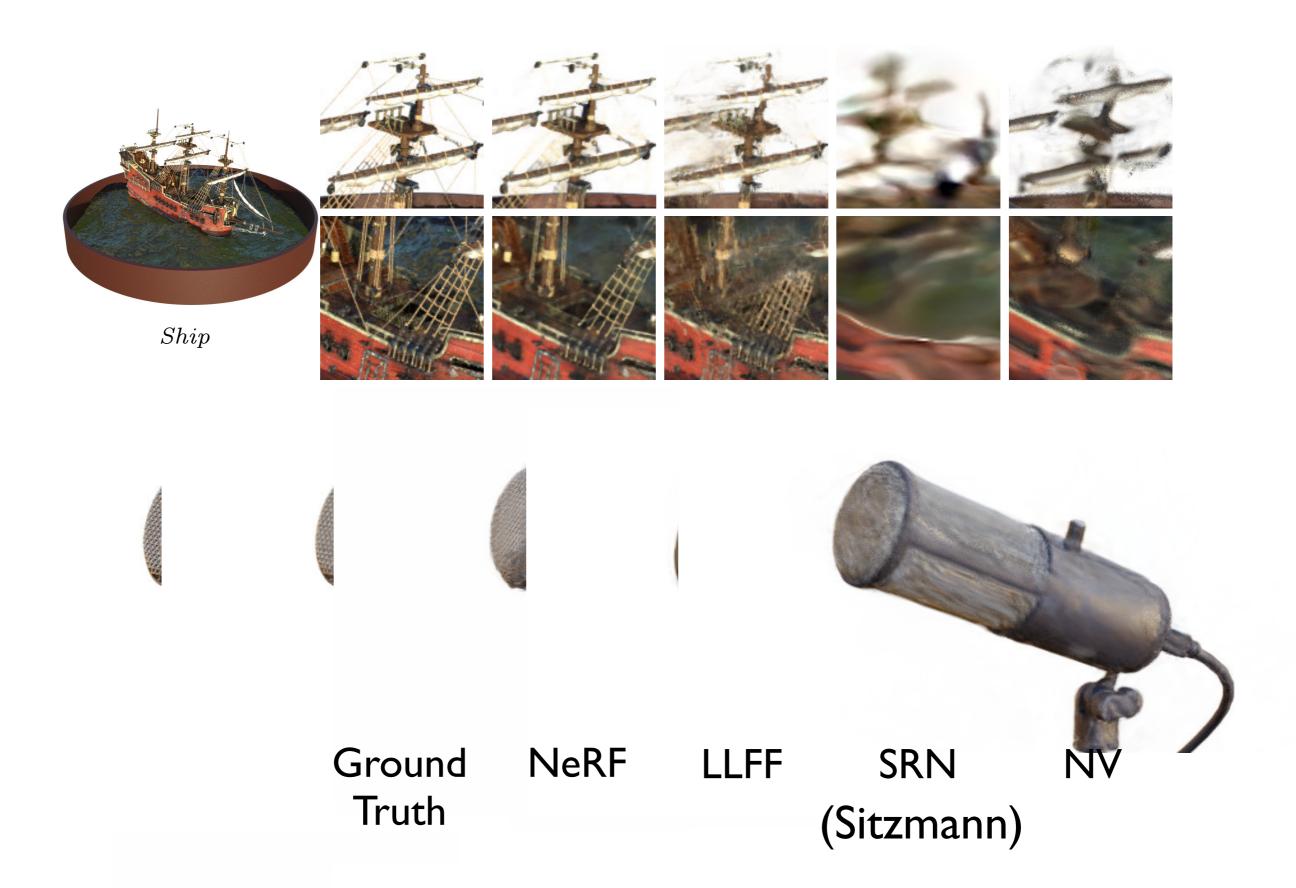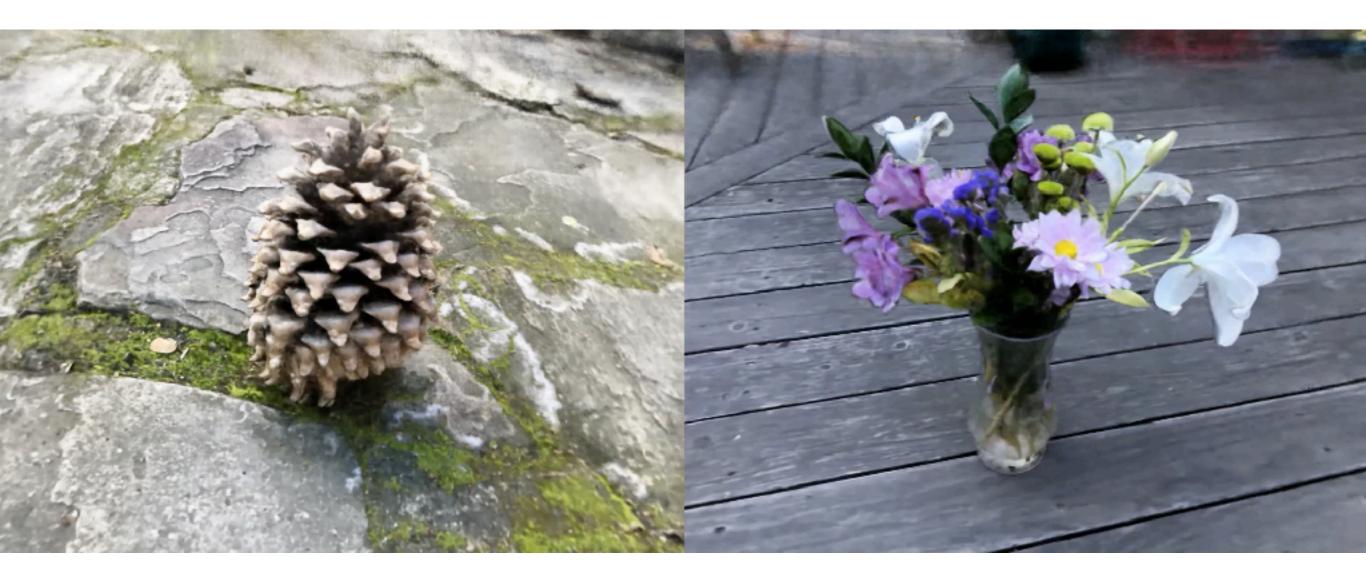- Another continuous scene representation using a FCN



$(x, y, z, \theta, \phi) \rightarrow F_{\Theta} \rightarrow (RGB\sigma)$

Predict density at each location, integrate along ray to get color (volume rendering)

[ NeRF, Mildenhall, Srinivasan, Tancik et al. 2020 ]

# Results



Ground Truth    NeRF    LLFF    SRN (Sitzmann)    NV

*Ship*

# Neural Radiance Fields

- Neural Radiance Fields, ~10s of input views



matthewtancik.com/nerf

# Next Lecture

- Image Generation, Generative Adversarial Networks