

ROB 498/599: Deep Learning for Robot Perception (DeepRob)

Lecture 20: 3D Vision Techniques; Nerf

03/31/2025



<https://deeprob.org/w25/>

Today

- Feedback and Recap (5min)
- 3D Shape Representations (40min)
- NeRF (25min)
- Summary and Takeaways (5min)

Tomorrow

Tuesday April 1, 2025

Starts 3:30pm EST @CSRB 2246

5min “lightning talk” - final project proposal presentations

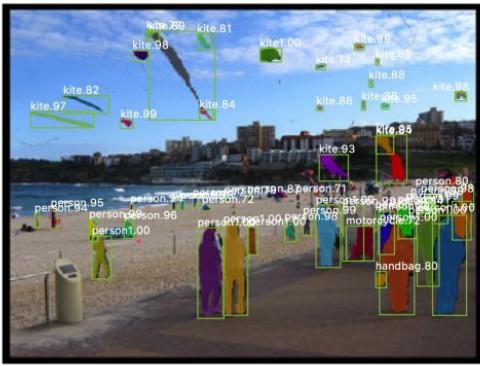
- All group members must speak
- Upload presentation files to [Final Project LightningTalk folder](#) on Google Drive by April 1st 3:30pm EST
- We will go by group number order
- 5% grade

3D Shape Representations

Previously... Predicting 3D Shapes of Objects

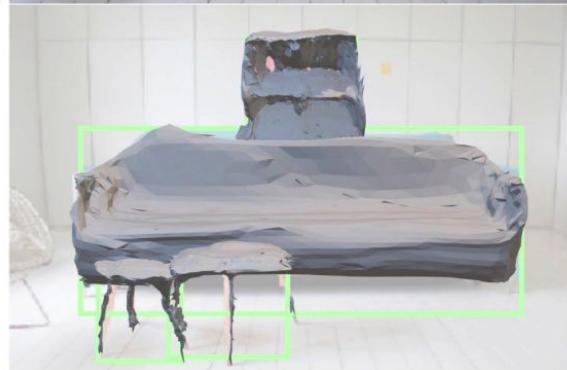
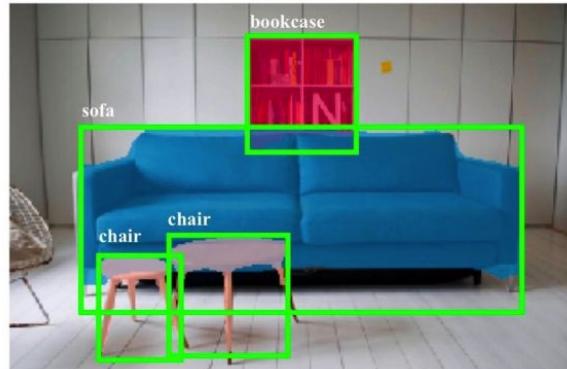
Mask R-CNN:

2D Image -> 2D shapes



Mesh R-CNN:

2D Image -> 3D shapes



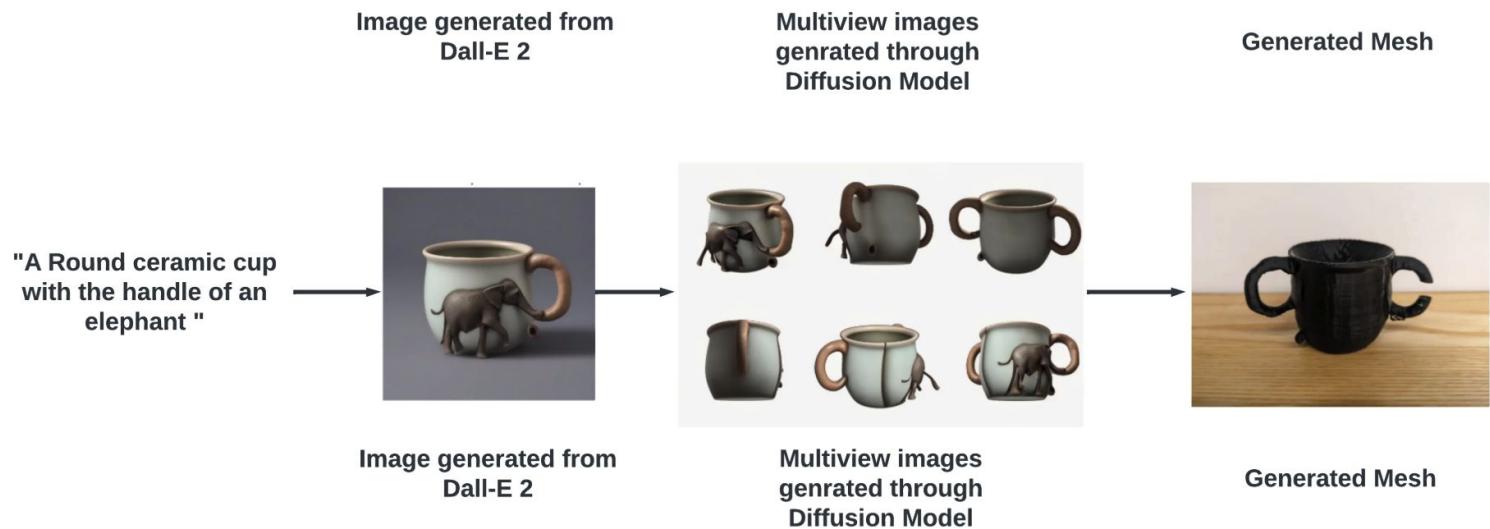
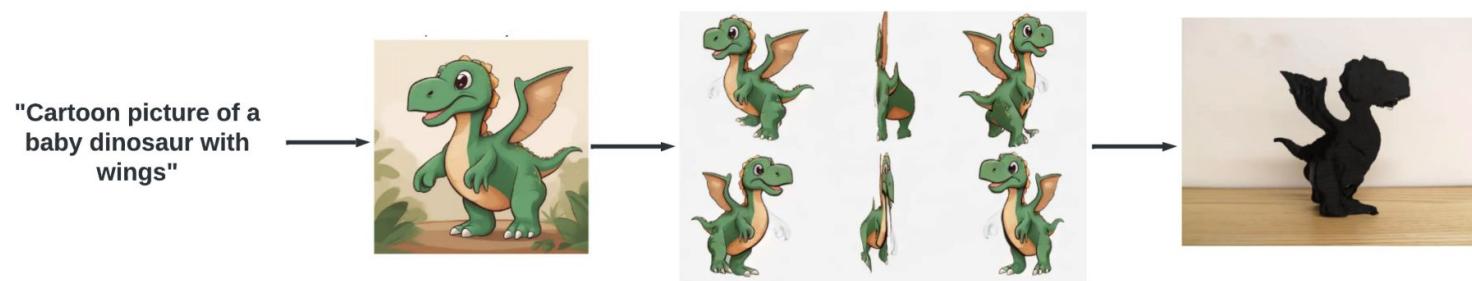
He, Gkioxari, Dollár, and Girshick, "Mask R-CNN", ICCV 2017

<https://arxiv.org/abs/1703.06870>

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

https://openaccess.thecvf.com/content_ICCV_2019/papers/Gkioxari_Mesh_R-CNN_ICCV_2019_paper.pdf

Also recall... “Speech2Mesh”



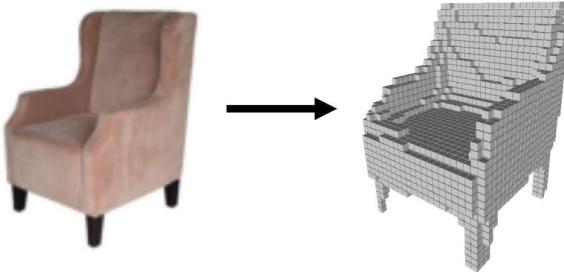
<https://deeprob.org/w24/reports/speech2mesh/>

DeepRob W24 Project

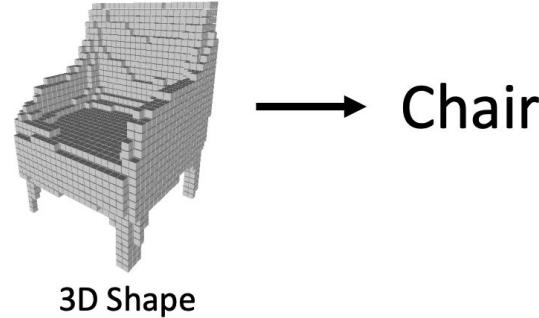
(Adeshara et al.)

Today

Predicting 3D Shapes
from single image



Processing 3D
input data

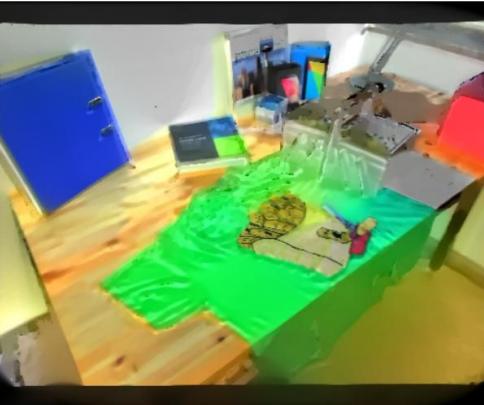


More topics in 3DV:

- Computing correspondences
- Multi-view stereo Structure from Motion
- Simultaneous Localization and Mapping (SLAM)
- Self-supervised learning
- Differentiable graphics
- 3D Sensors

<https://3dvconf.github.io/2025/>
(3DV)
<https://s2025.siggraph.org/>
(SIGGRAPH)

3D data modalities



<https://www.nist.gov/image/hyperspectralcube4501png>

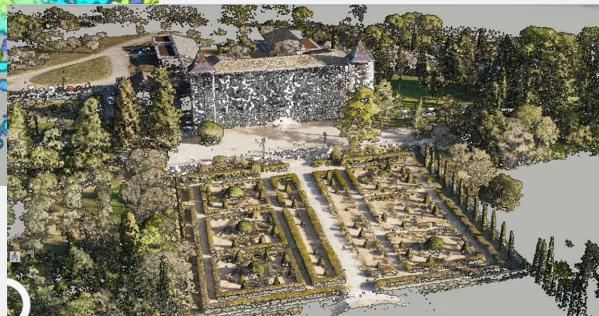
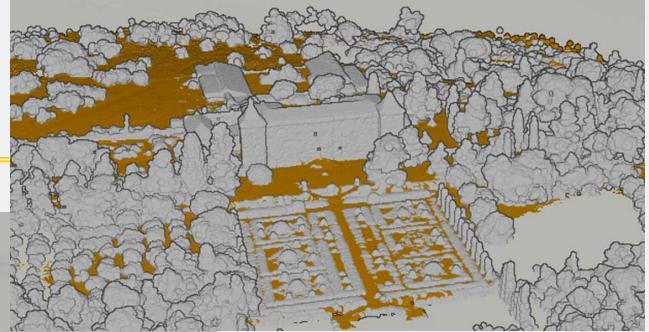
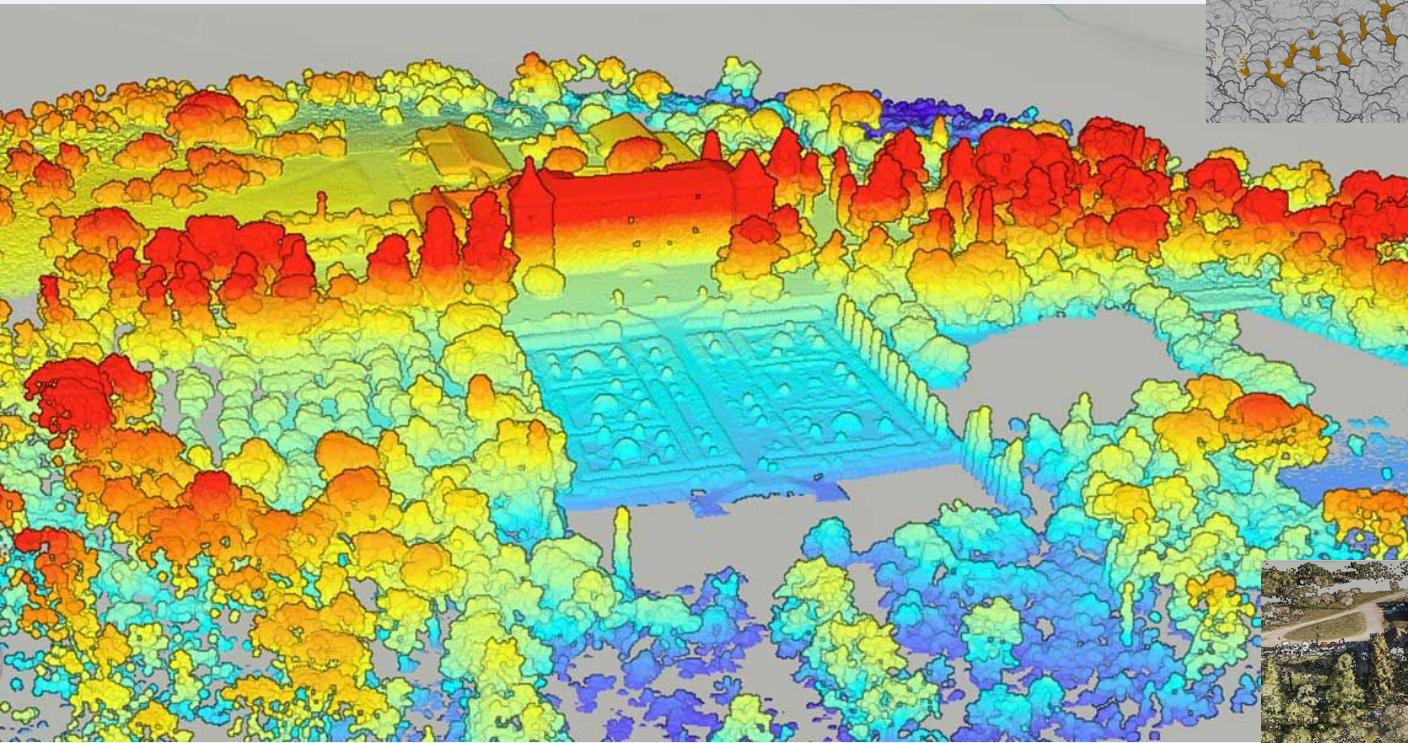
<https://www.flir.com/oem/adas/adas-data-set-form/?srsltid=AfmBOooTeV1G2vPoiodzG9ciTWSL18DTQSWGiuV5qXccWXR2g39S>



RGB-T
(thermal)

<https://paperswithcode.com/dataset/tum-rgb-d>

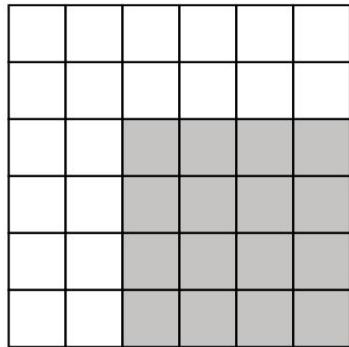
3D data modalities



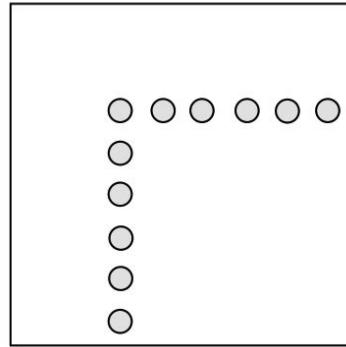
<https://www.yellowscan.com/knowledge/lidar-point-cloud-basics/>

3D Shape Representations

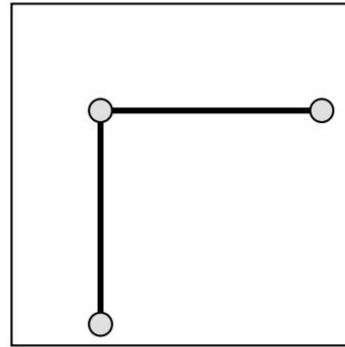
8
8
2
2
2
2



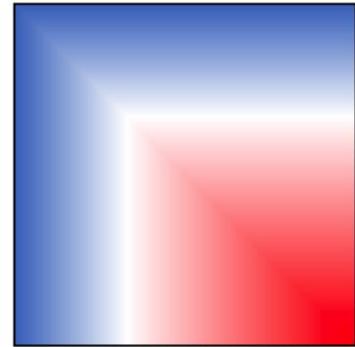
Depth
Map



Pointcloud



Mesh



Implicit
Surface

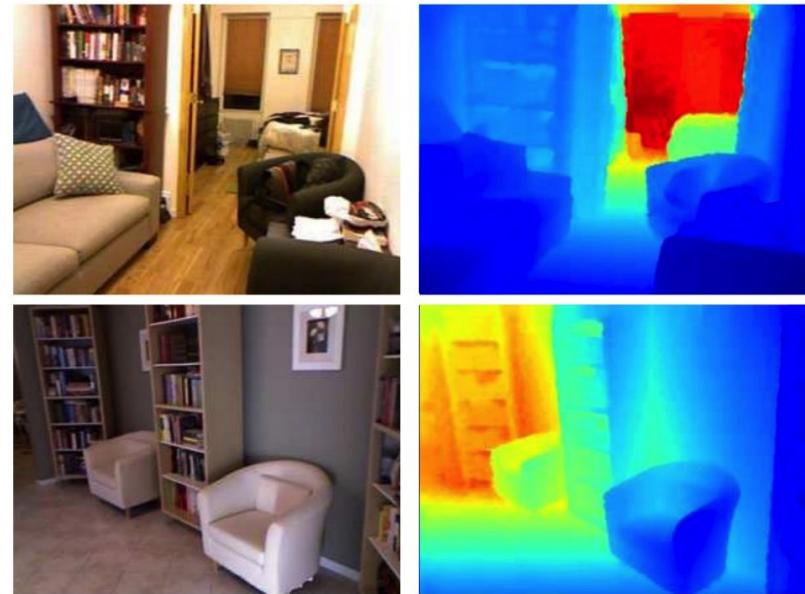
3D Shape Representations

Q: Why 2.5D?

For each pixel, **depth map** gives distance from the camera to the object in the world at that pixel

RGB image + Depth image
= RGB-D Image **(2.5D)**

This type of data can be recorded directly for some types of 3D sensors (e.g. Microsoft Kinect)

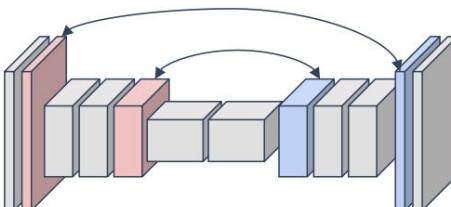


RGB Image: $3 \times H \times W$ Depth Map: $H \times W$

https://openaccess.thecvf.com/content_iccv_2015/papers/Eigen_Predicting_Depth_Surface_ICCV_2015_paper.pdf

Predicting Depth Maps

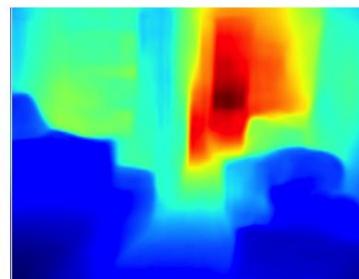
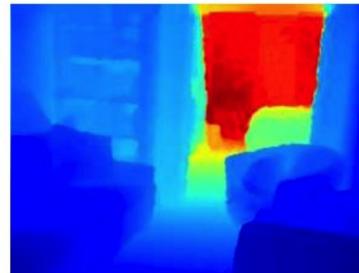
(KITTI dataset)



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

**Fully Convolutional
network**

Predicted Depth Image:
 $1 \times H \times W$



Predicted Depth Image:
 $1 \times H \times W$

https://papers.nips.cc/paper_2014/hash/91c56ce4a249fae5419b90cba831e303-Abstract.html

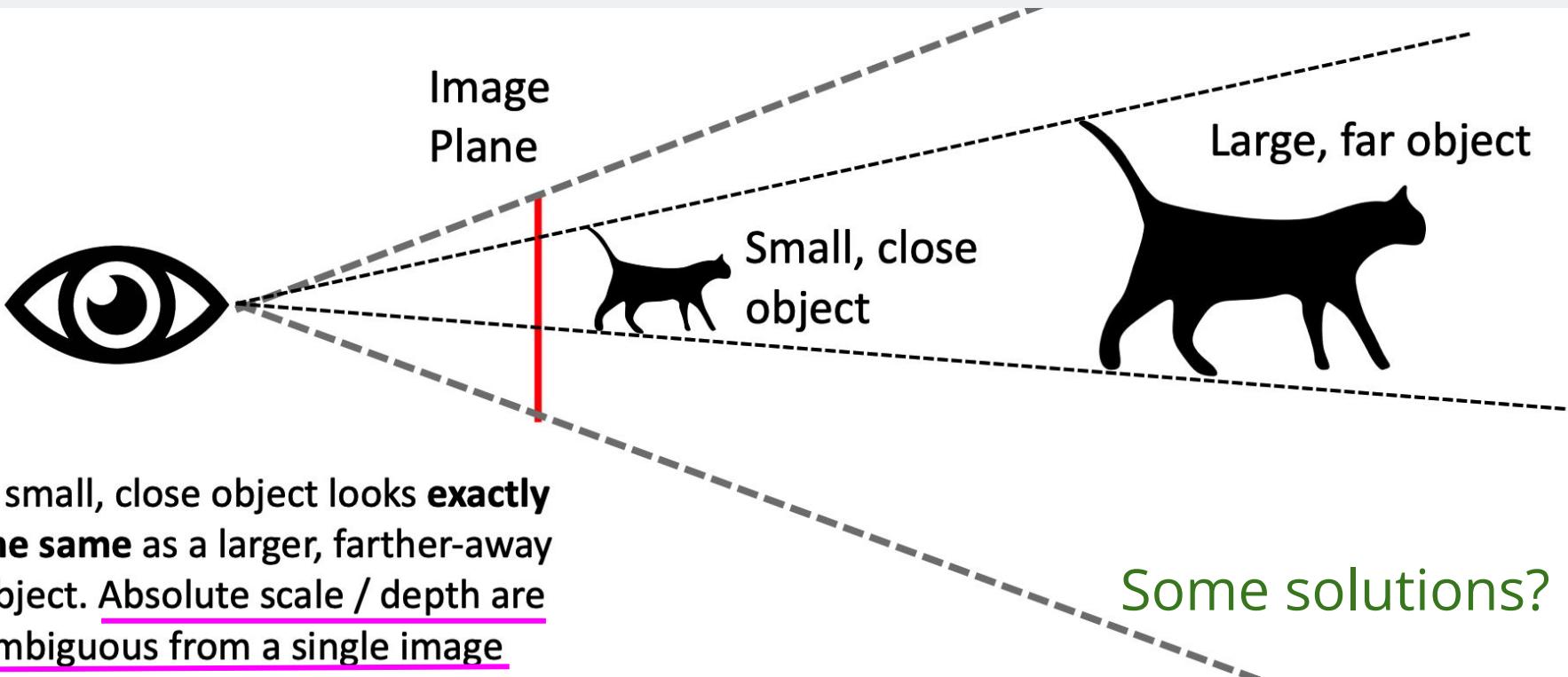
<https://www.cvlibs.net/datasets/kitti/>

Aha Slides (In-class participation)

<https://ahaslides.com/YSRO4>



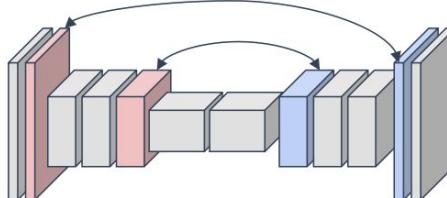
Problem: Scale / Depth Ambiguity



Predicting Depth Maps

Scale-invariant mean squared error
(in log space)

$$\begin{aligned} D(y, y^*) &= \frac{1}{2n^2} \sum_{i,j} ((\log y_i - \log y_j) - (\log y_i^* - \log y_j^*))^2 \\ &= \frac{1}{n} \sum_i d_i^2 - \frac{1}{n^2} \sum_{i,j} d_i d_j = \frac{1}{n} \sum_i d_i^2 - \frac{1}{n^2} \left(\sum_i d_i \right)^2 \end{aligned}$$

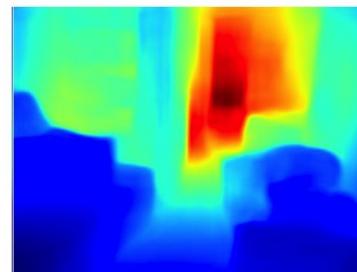
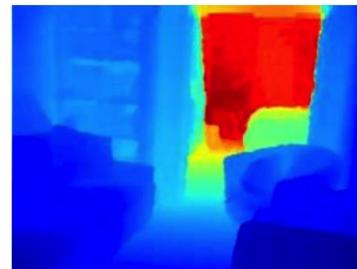


RGB Input Image:
3 x H x W

Fully Convolutional
network

Predicted Depth Image:

1 x H x W



Per-Pixel Loss
(Scale invariant)

Predicted Depth Image:
1 x H x W

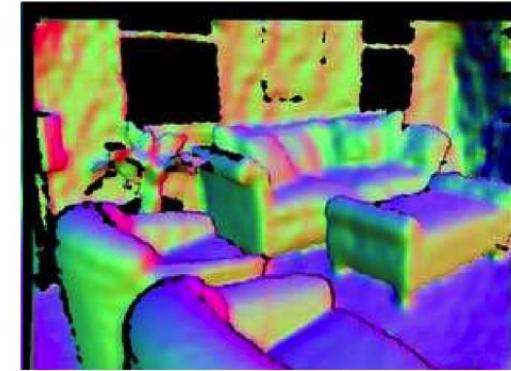
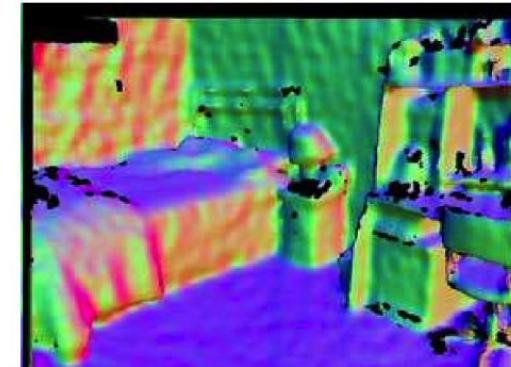
3D Shape Representations: Surface Normals

For each pixel, **surface normals** give a vector giving the normal vector to the object in the world for that pixel

Q: How to obtain surface normals?



RGB Image: $3 \times H \times W$



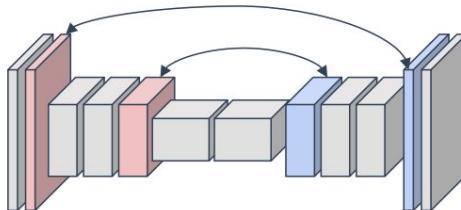
Normals: $3 \times H \times W$

<https://arxiv.org/pdf/1411.4734>

Predicting Surface Normals

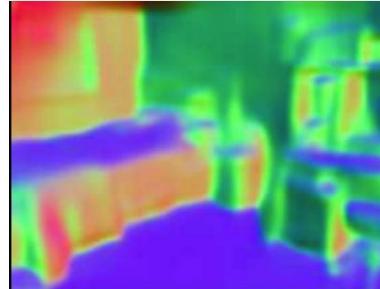
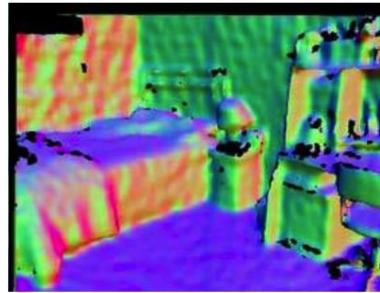


RGB Input Image:
 $3 \times H \times W$



**Fully Convolutional
network**

Ground-truth Normals:
 $3 \times H \times W$



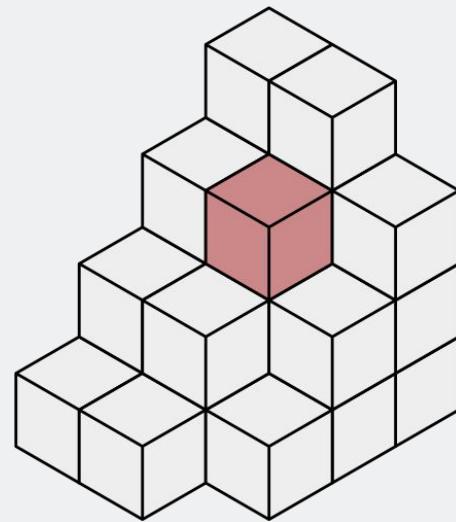
Predicted Normals:
 $3 \times H \times W$

Per-Pixel Loss:
 $(x \cdot y) / (|x| |y|)$

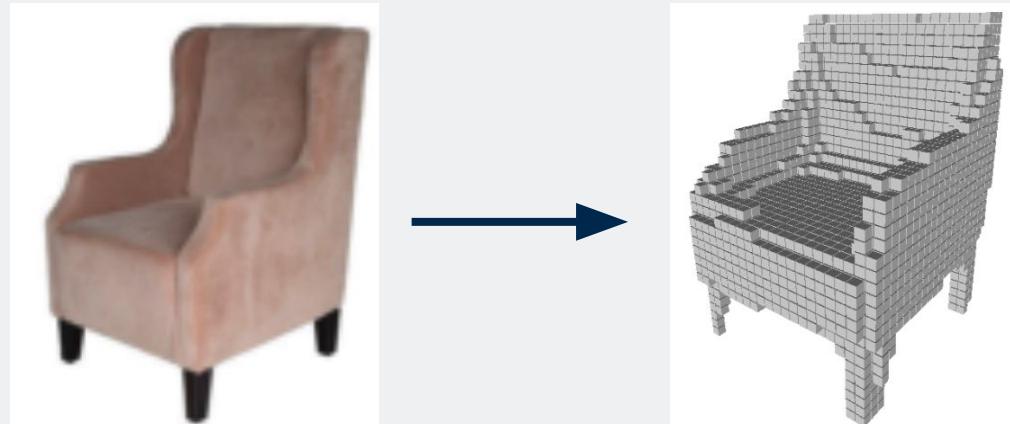
Recall:

$$\begin{aligned} x \cdot y \\ = |x| |y| \cos \theta \end{aligned}$$

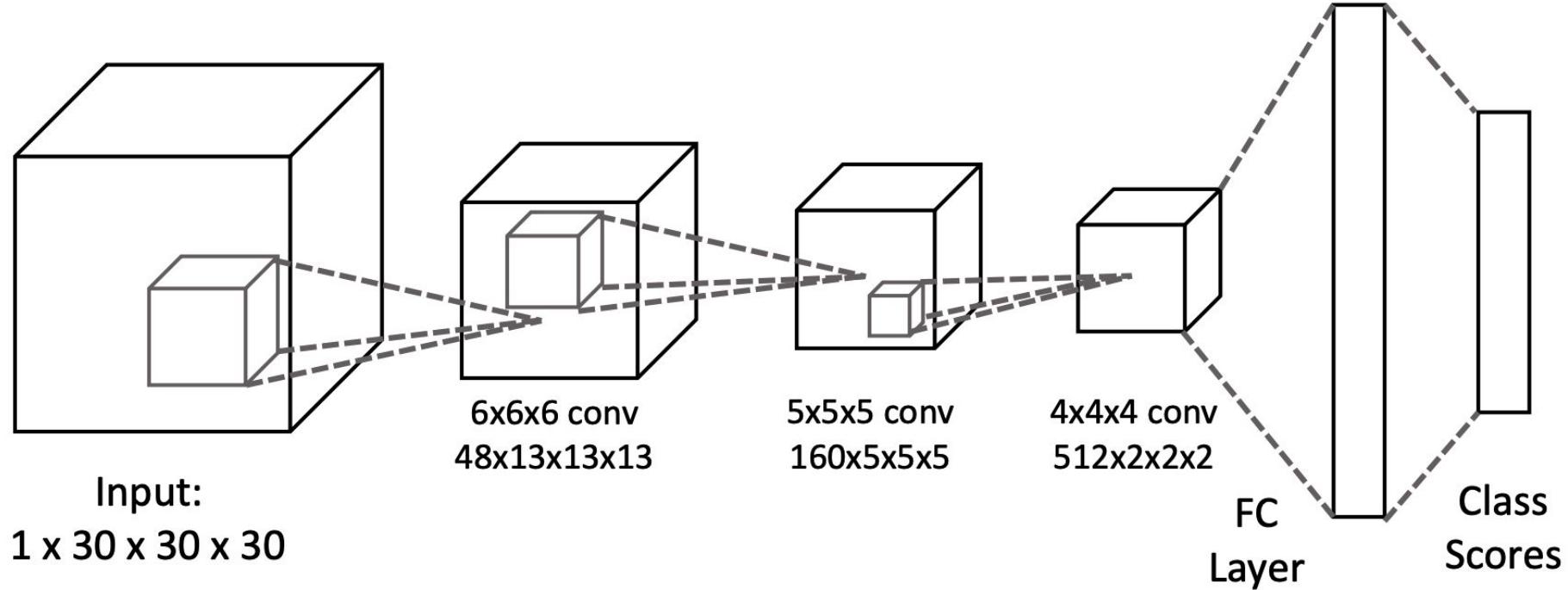
3D Shape Representations: Voxels



- Represent a shape with a $V \times V \times V$ 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!



Processing Voxel Inputs: 3D Convolution



Input:
 $1 \times 30 \times 30 \times 30$

$6 \times 6 \times 6$ conv
 $48 \times 13 \times 13 \times 13$

$5 \times 5 \times 5$ conv
 $160 \times 5 \times 5 \times 5$

$4 \times 4 \times 4$ conv
 $512 \times 2 \times 2 \times 2$

FC
Layer

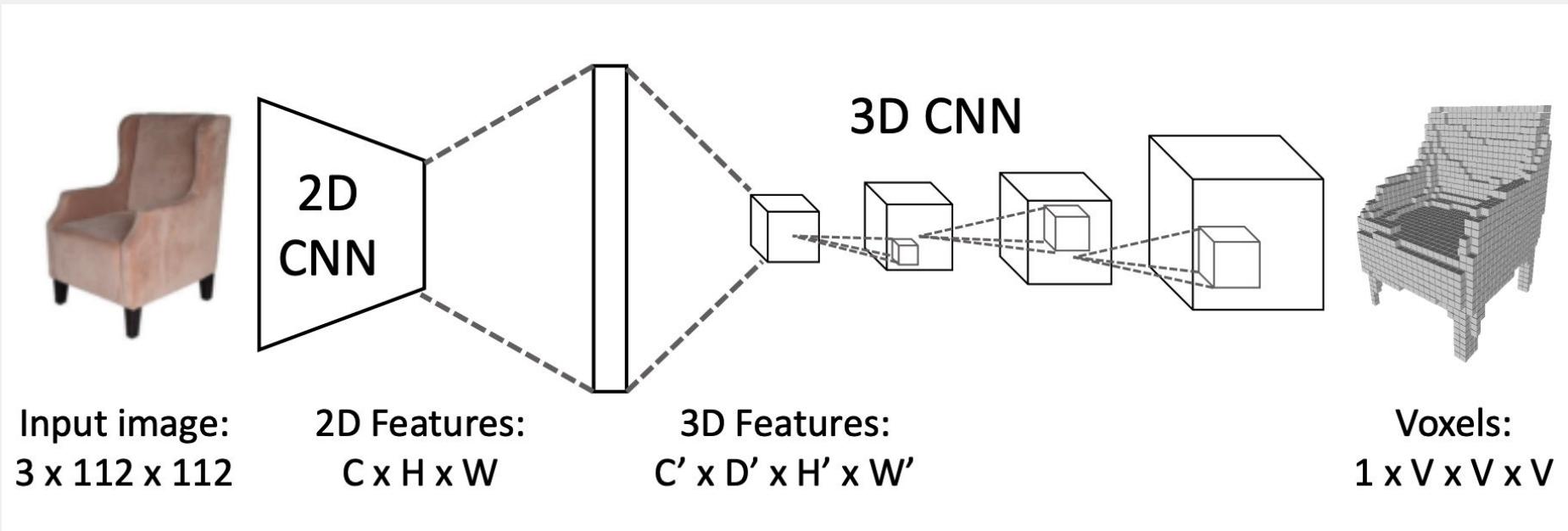
Class
Scores

Train with classification loss

<https://arxiv.org/pdf/1406.5670>

<https://3dshapenets.cs.princeton.edu/>

Generating Voxel Shapes: 3D Convolution



<https://arxiv.org/pdf/1604.00449>

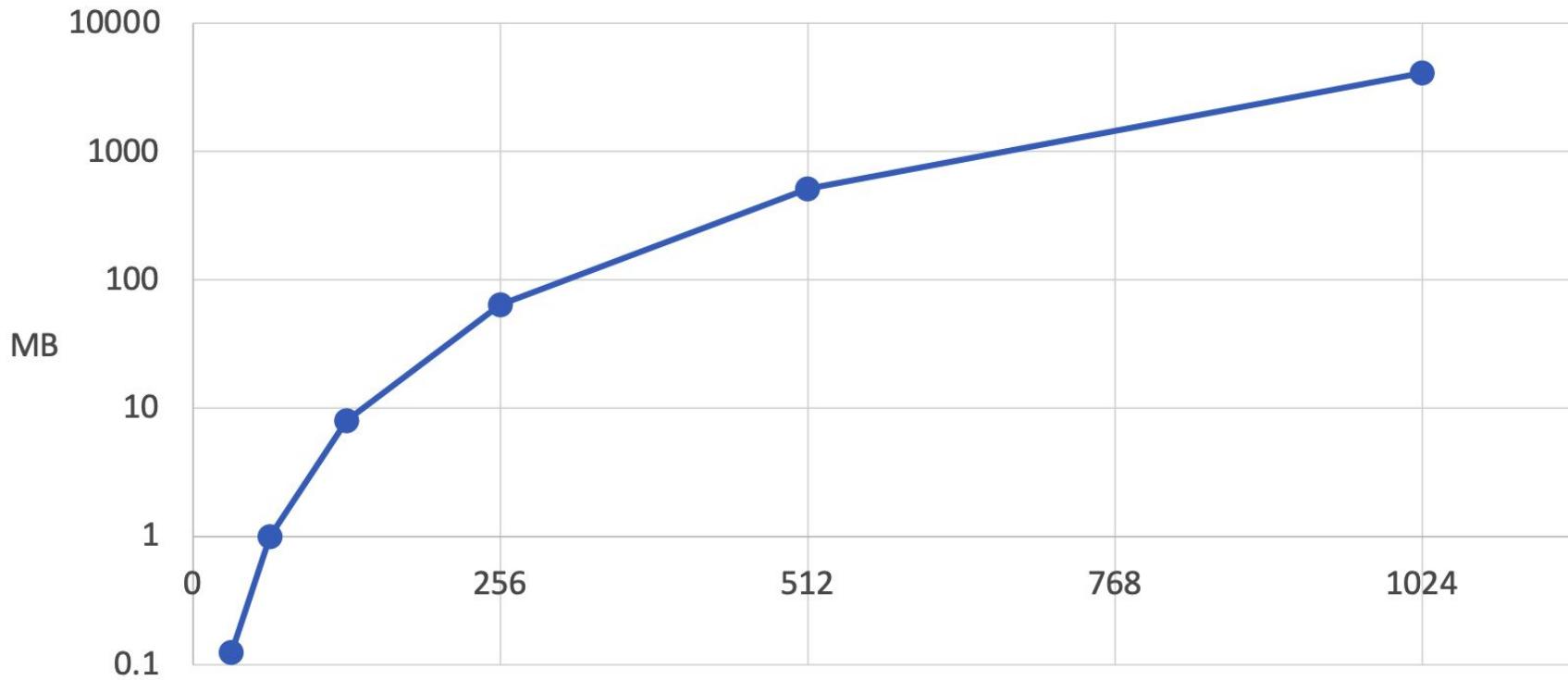
Train with per-voxel cross-entropy loss

Q: voxels with occluded parts?

Voxel: Memory Usage

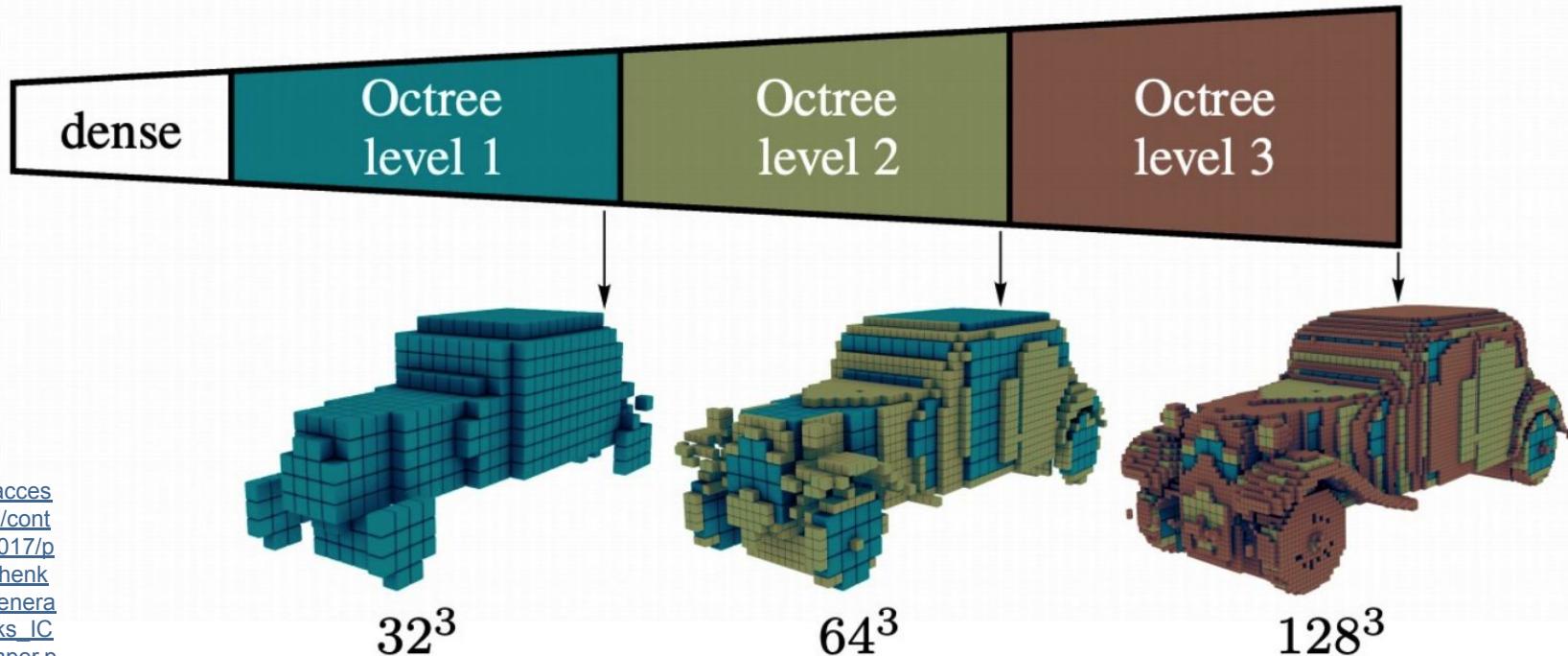
Storing 1024^3 voxel grid takes 4GB of memory!

Voxel memory usage ($V \times V \times V$ float32 numbers)



Scaling Voxels: Oct-Trees

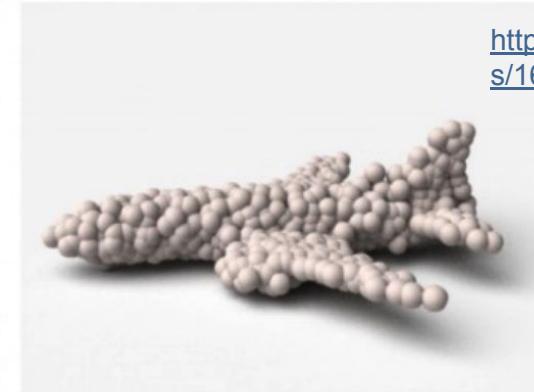
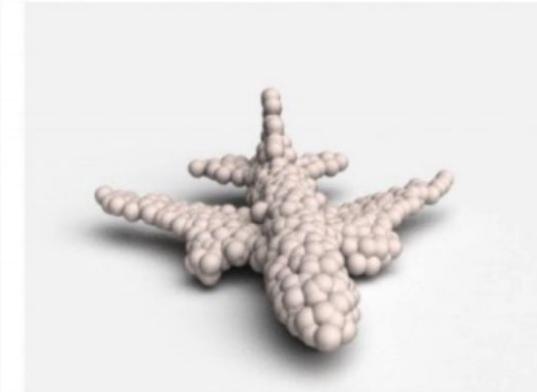
Use voxel grids with heterogenous resolution!



https://openaccess.thecvf.com/content_ICCV_2017/papers/Tatarchenko_Octree_Generating_Networks_ICCV_2017_paper.pdf

3D Shape Representations: Point Cloud

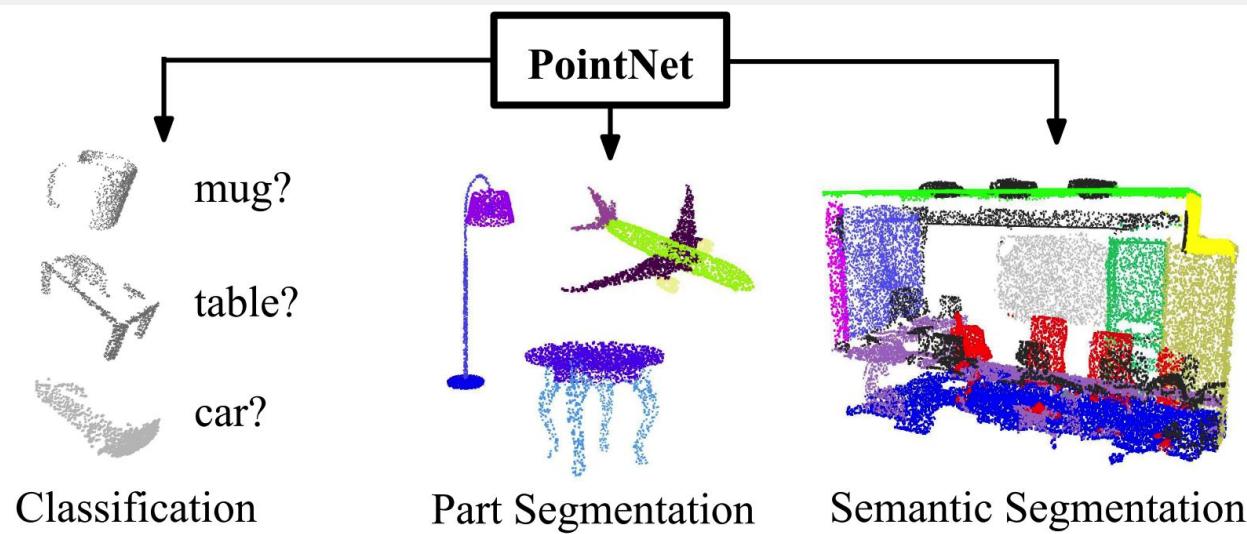
- Represent shape as a set of P points in 3D space
- (+) Can represent fine structures without huge numbers of points
- (-) Requires new architecture, losses, etc
- (-) Doesn't explicitly represent the surface of the shape: extracting a mesh for rendering or other applications requires post-processing



<https://arxiv.org/abs/1612.00603>



PointNet



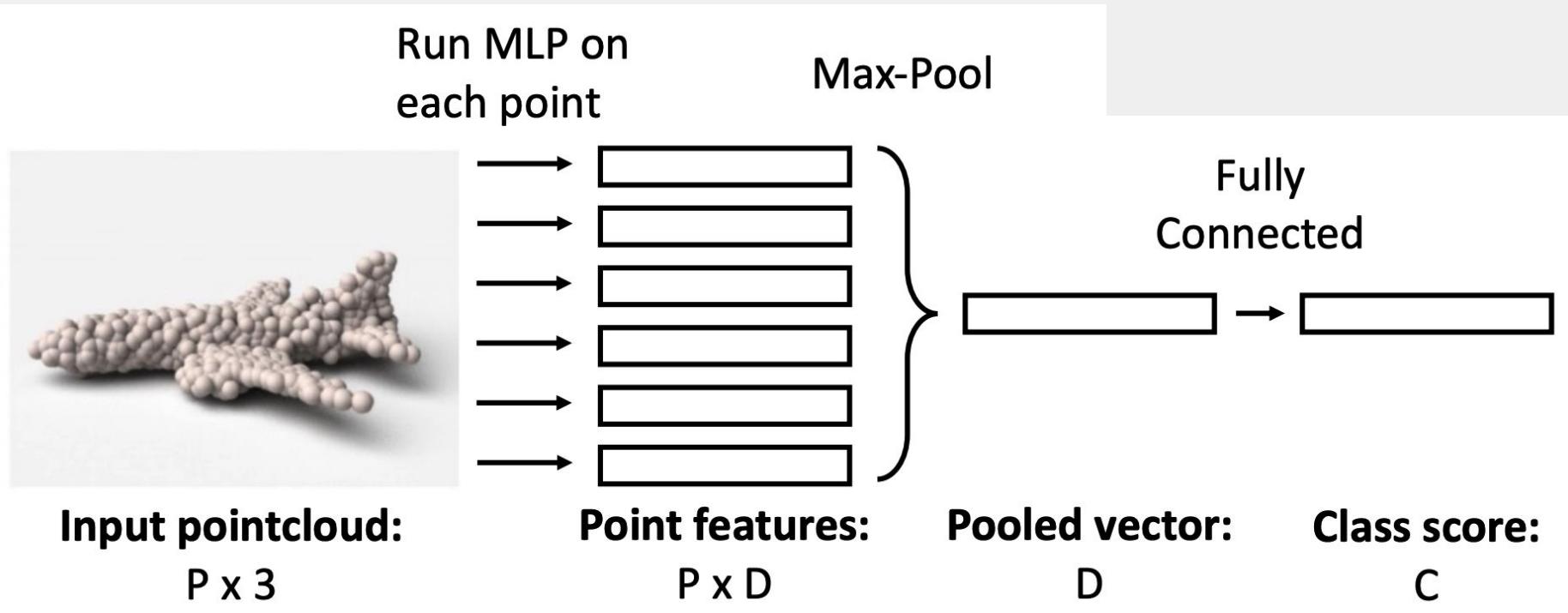
Classification Part Segmentation Semantic Segmentation

Figure 1. Applications of PointNet. We propose a novel deep net architecture that consumes raw point cloud (set of points) without voxelization or rendering. It is a unified architecture that learns both global and local point features, providing a simple, efficient and effective approach for a number of 3D recognition tasks.

https://openaccess.thecvf.com/content_cvpr_2017/papers/Qi_PointNet_Deep_Learning_CVPR_2017_paper.pdf

PointNet

Q: does order matter?



PointNet and more

PointNet (3D point cloud processing, no voxel or rendering):

https://openaccess.thecvf.com/content_cvpr_2017/papers/Qi_PointNet_Deep_Learning_CVPR_2017_paper.pdf

PointNet++ (adaptively combine features from multiple scales):

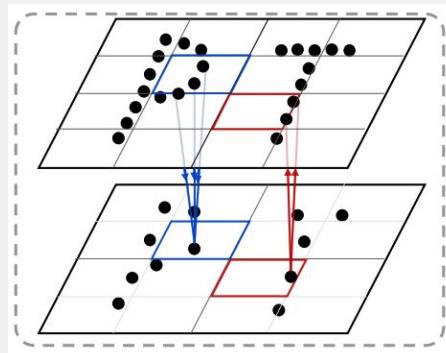
<https://arxiv.org/pdf/1706.02413>

PointNeXt (Improved Training and Scaling Strategies):

https://proceedings.neurips.cc/paper_files/paper/2022/file/9318763d049edf9a1f2779b2a59911d3-Paper-Conference.pdf

Point Transformer (v2) (partition-based group vector attention):

https://proceedings.neurips.cc/paper_files/paper/2022/file/d78ece6613953f46501b958b7bb4582f-Paper-Conference.pdf



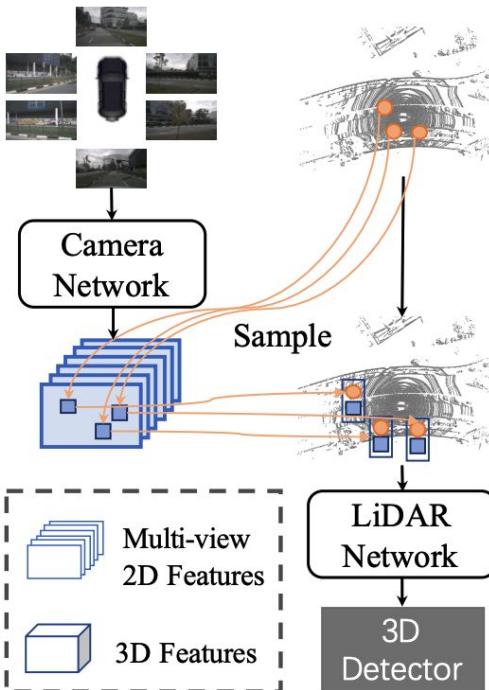
PointNet and more

Applications (e.g., LiDAR+Camera Fusion)

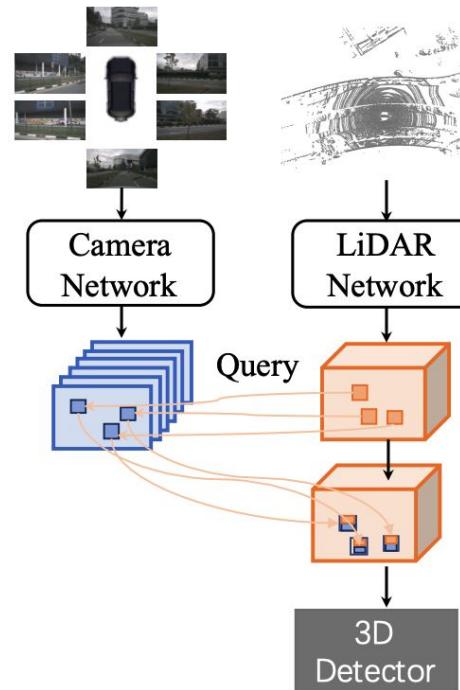
BEVFusion

https://proceedings.neurips.cc/paper_files/paper/2022/file/43d2b7fbee8431f7cef0d0afed51c691-Paper-Conference.pdf

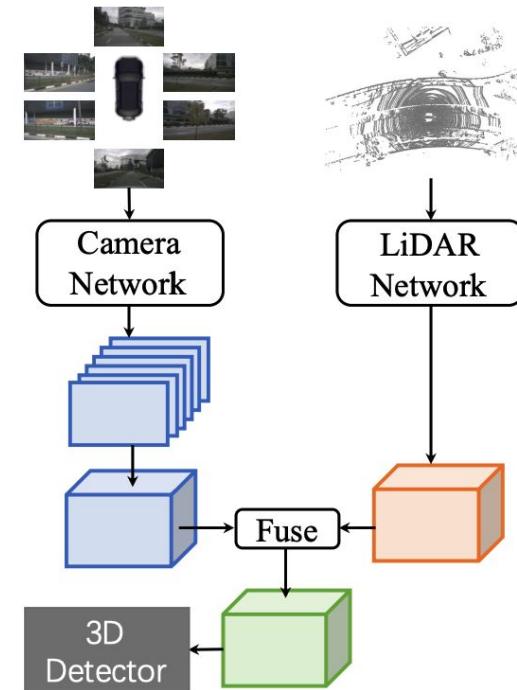
(a) Point-level Fusion



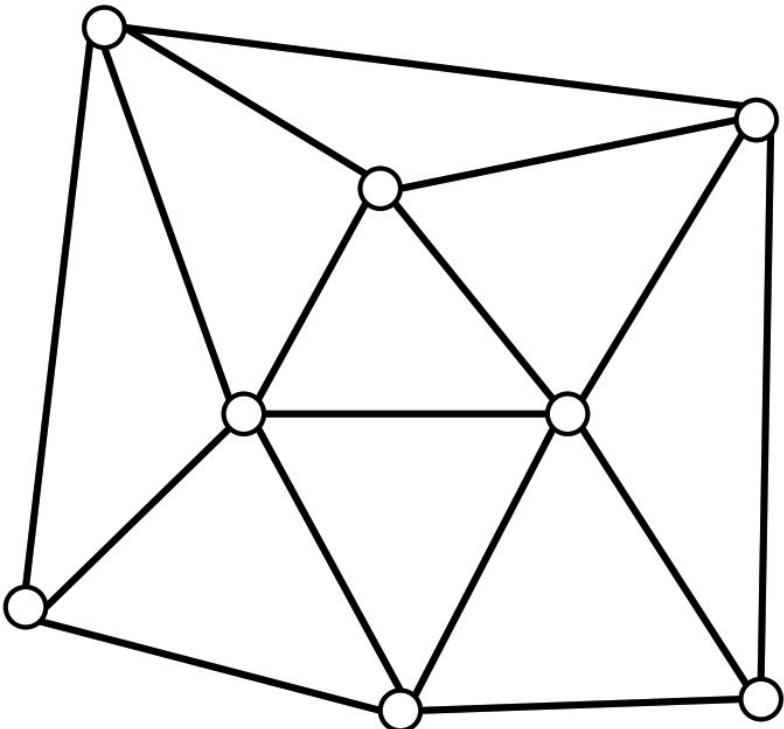
(b) Feature-level Fusion



(c) Our BEVFusion



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

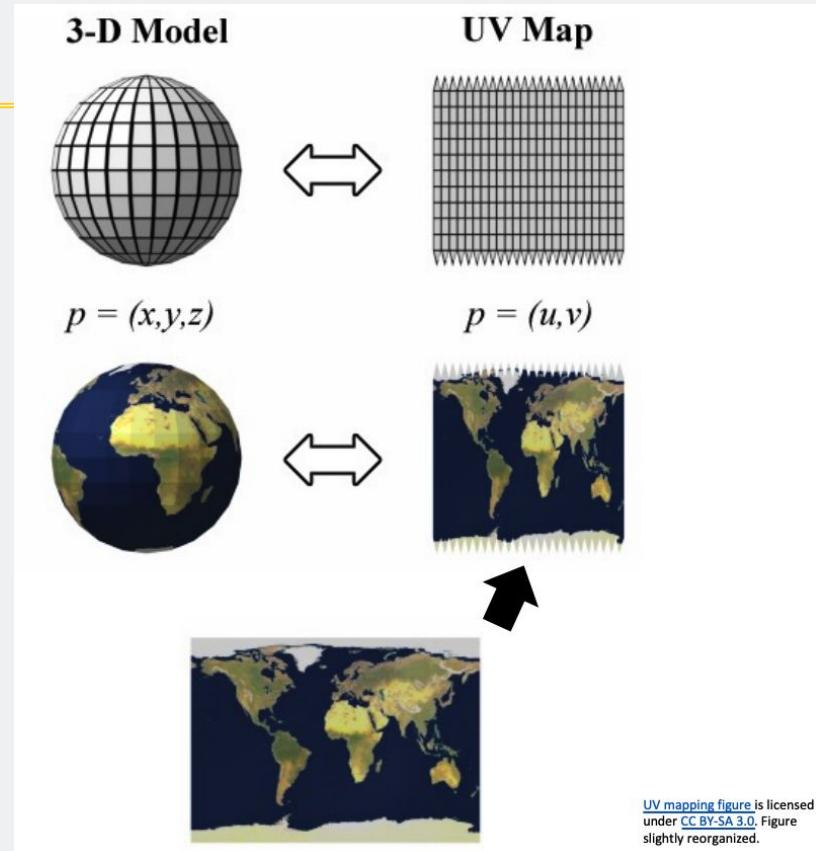
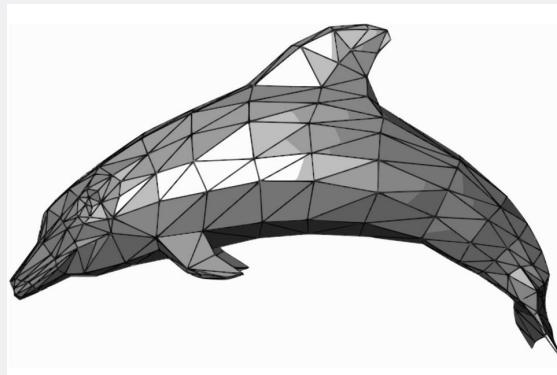
(+) Adaptive: Can represent flat surfaces very efficiently, can allocate more faces to areas with fine detail

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

Examples

<https://graphics.stanford.edu/data/3Dscanrep/>

“Stanford
bunny”



SMPL

<https://smpl.is.tue.mpg.de/>

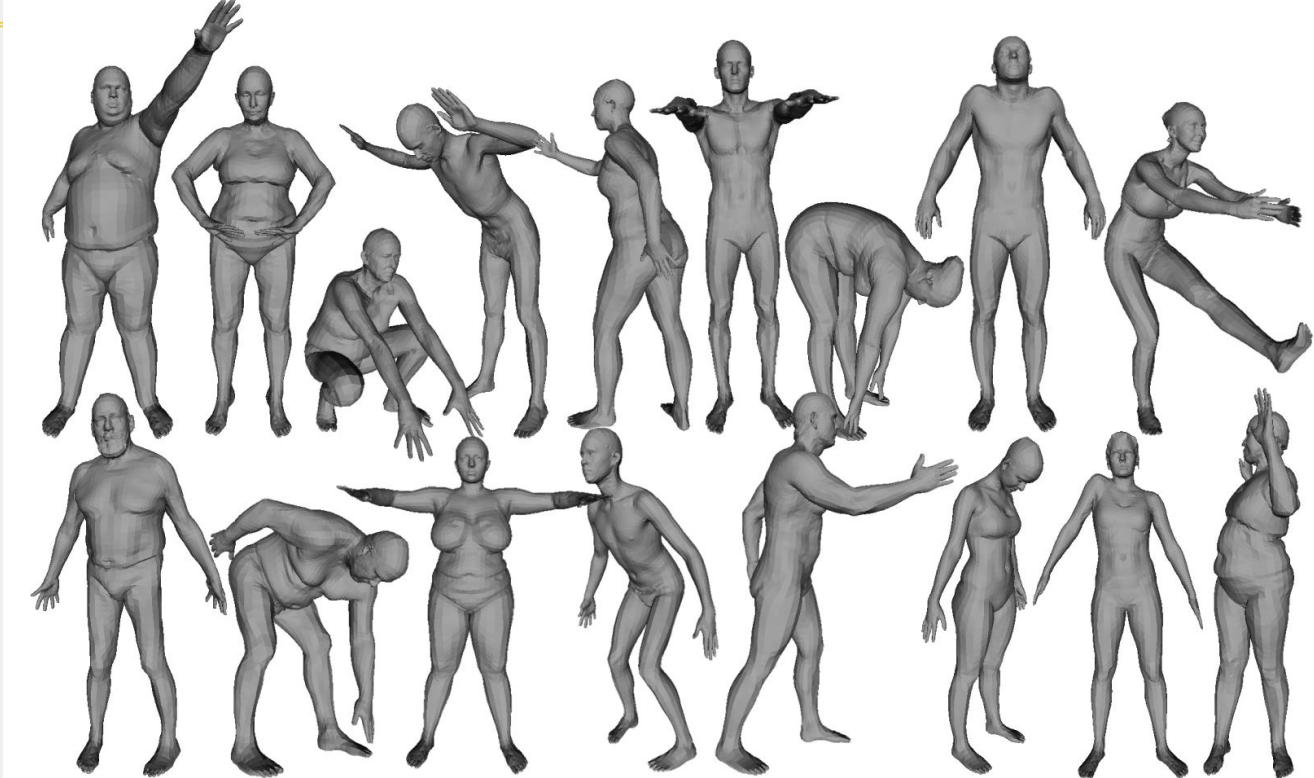
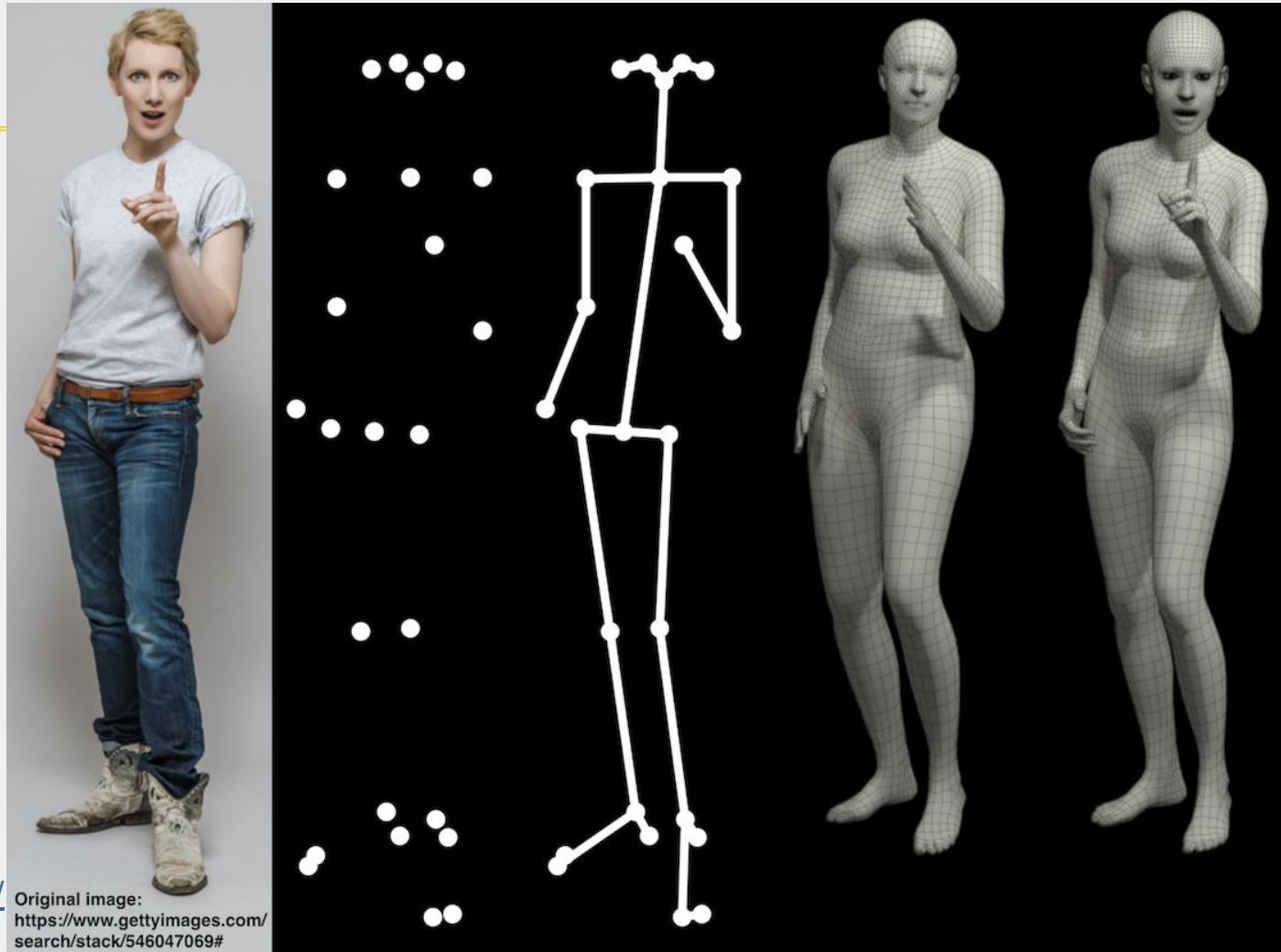


Figure 4: Sample registrations from the multipose dataset.

SMPL-X



<https://smpl-x.is.tue.mpg.de/>

Original image:
[https://www.gettyimages.com/
search/stack/546047069#](https://www.gettyimages.com/search/stack/546047069#)

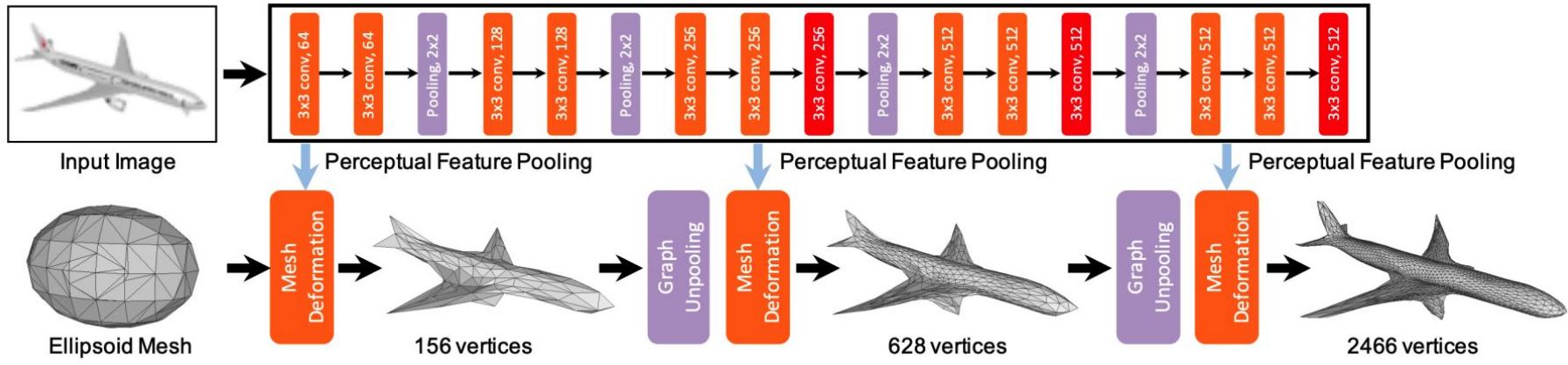
Predicting Meshes: Pixel2Mesh

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



<https://arxiv.org/pdf/1804.01654>

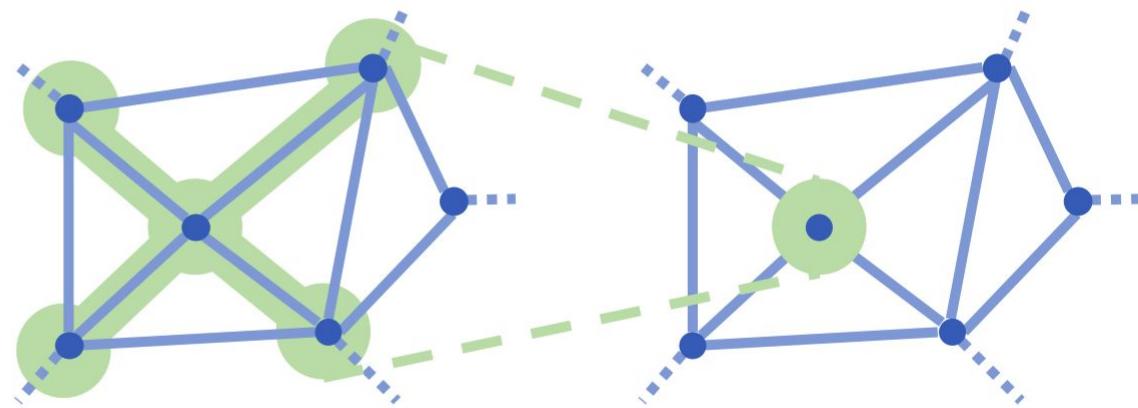
Graph Convolution

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

Vertex v_i has feature f_i

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

Use same weights W_0 and W_1 to compute all outputs



Input: Graph with a feature vector at each vertex

Output: New feature vector for each vertex

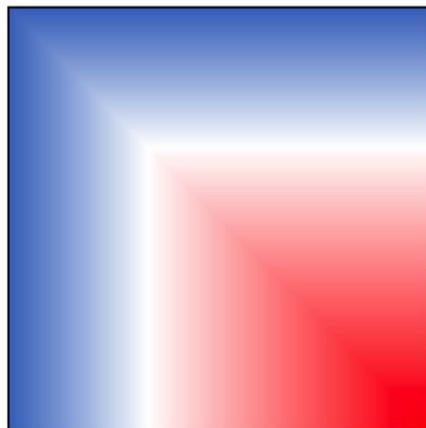
3D Shape Representations: Implicit Functions

Learn a function to classify arbitrary 3D points as inside / outside the shape

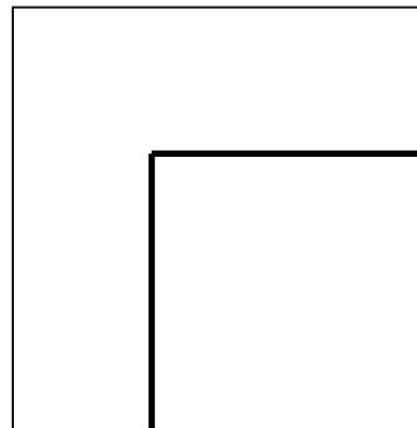
$$o : \mathbb{R}^3 \rightarrow \{0, 1\}$$

The surface of the 3D object is the level set

$$\{x : o(x) = \frac{1}{2}\}$$



Implicit function



Explicit Shape

Binary classification problem

- Occupancy
- SDF (signed distance function)

Neural Radiance Field (NeRF) For View Synthesis

<https://arxiv.org/pdf/2003.08934>

View Synthesis

<https://www.matthewtancik.com/nerf>

Input: Many images of the same scene
(with known camera parameters)

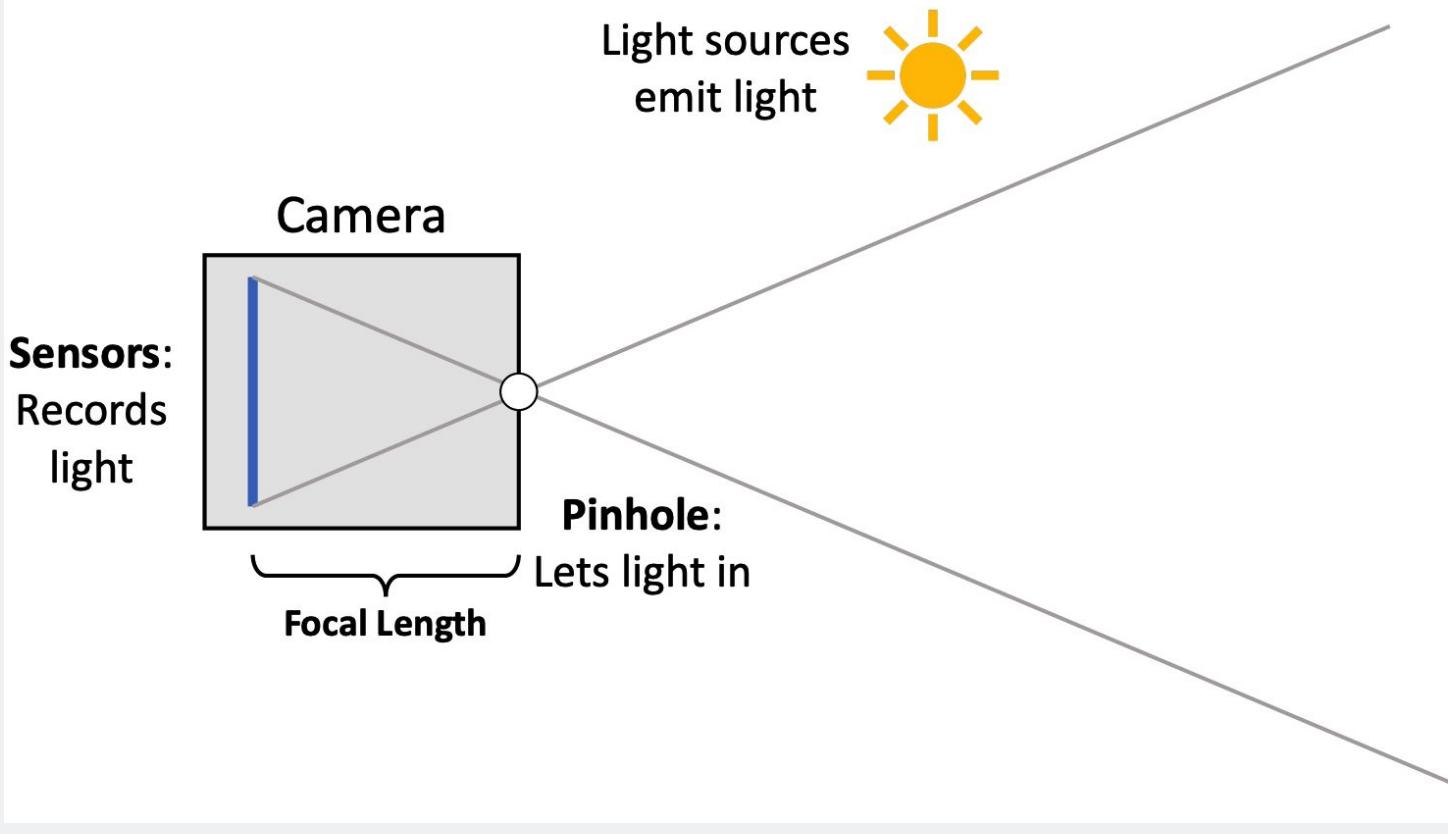


Output: Images showing the scene from novel viewpoints



Scene rendering

Stepping Back: Pinhole Camera Model

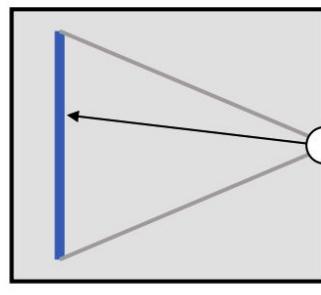


Stepping Back: Pinhole Camera Model

Reflected light
captured by
the camera

Camera

Sensors:
Records
light



Pinhole:
Lets light in

Light sources
emit light



Scenario 1

Objects in the
world reflect light

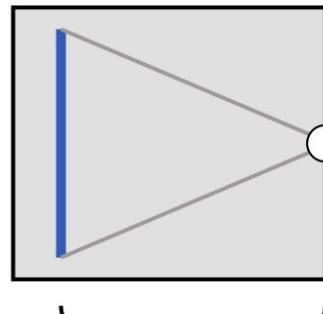


Stepping Back: Pinhole Camera Model

Reflected light
captured by
the camera

Camera

Sensors:
Records
light



Pinhole:
Lets light in

Light sources
emit light



Opaque objects
block light



Scenario 2

Objects in the
world reflect light

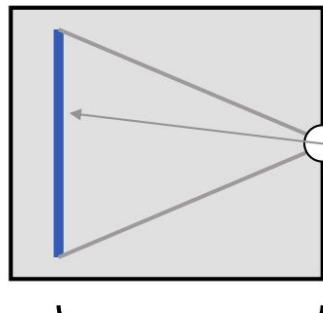


Stepping Back: Pinhole Camera Model

Reflected light
captured by
the camera

Camera

Sensors:
Records
light

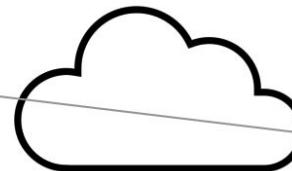


Pinhole:
Lets light in

Light sources
emit light



Transparent objects
attenuate light



Scenario 3

Objects in the
world reflect light

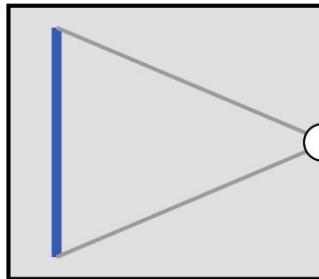


Volume Rendering

Abstract away light sources, objects.

For each point in space, need to know:

- (1) How much light does it emit?
- (2) How opaque is it? $\sigma \in [0,1]$

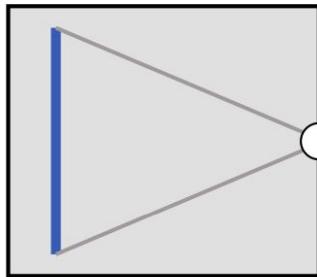


Stepping Back: Pinhole Camera Model

Abstract away light sources, objects.

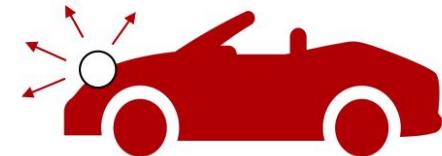
For each point in space, need to know:

- (1) How much light does it emit?
- (2) How opaque is it? $\sigma \in [0,1]$

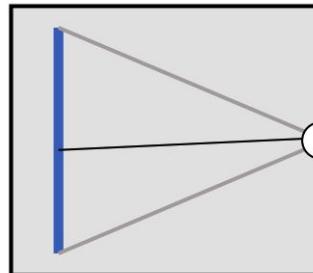


Point in empty space:
(1) Emits no light (black)
(2) Completely transparent $\sigma = 0$

Point on car:
(1) Emits red light in hemisphere
(2) Complete opaque
 $\sigma = 1$



Volume Rendering



Ray origin

Parameterize each ray as origin plus

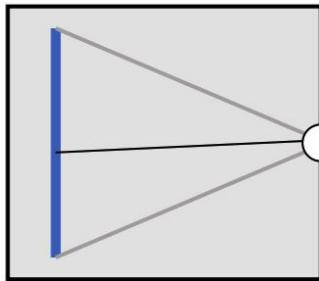
direction: $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$

Volume Density is $\sigma(\mathbf{p}) \in [0,1]$

Color that a point \mathbf{p} emits in direction \mathbf{d}

is $c(\mathbf{p}, \mathbf{d}) \in [0,1]^3$

Volume Rendering



Parameterize each ray as origin plus direction: $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$

Volume Density is $\sigma(\mathbf{p}) \in [0,1]$

Color that a point \mathbf{p} emits in direction \mathbf{d} is $c(\mathbf{p}, \mathbf{d}) \in [0,1]^3$

Color observed by the camera given by **volume rendering equation**:

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) c(\mathbf{r}(t), \mathbf{d}) dt$$

Ray origin

Near point: t_n

Current point: t

Far point: t_f

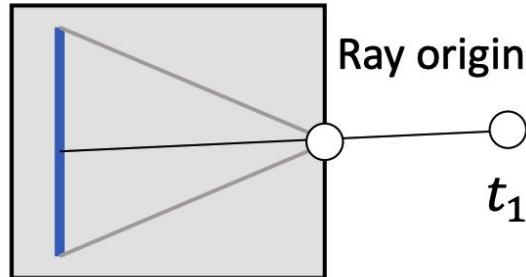
Transmittance: How much light from the current point will reach the camera?

Opacity: How opaque is the current point?

Color: What color does the current point emit along the direction toward the camera?

Transmittance

Compute transmittance by
**accumulating volume
density** up to current point



Parameterize each ray as origin plus
direction: $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$

Volume Density is $\sigma(\mathbf{p}) \in [0,1]$

Color that a point \mathbf{p} emits in direction \mathbf{d}
is $c(\mathbf{p}, \mathbf{d}) \in [0,1]^3$

Color observed by the camera given
by **volume rendering equation**:

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt$$

$$T(t) = \exp \left(- \int_{t_n}^t \sigma(\mathbf{r}(s)) ds \right)$$

Approximate integrals with a set of samples:

$$C(\mathbf{r}) \approx \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i$$

$$T_i = \exp \left(- \sum_{j=1}^{i-1} \sigma_j \delta_j \right)$$

Mildenhall et al, "Representing
Scenes as Neural Radiance Fields
for View Synthesis", ECCV 2020

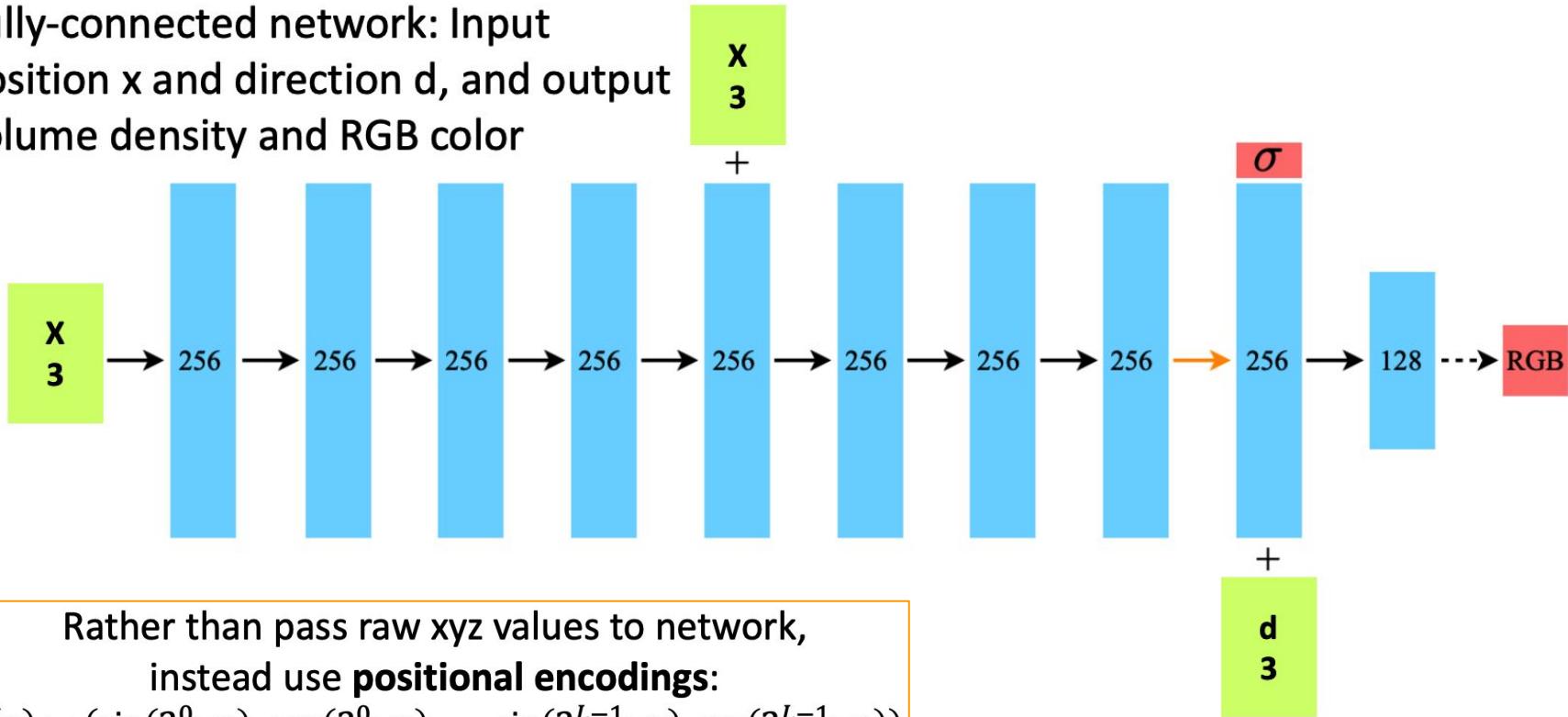
NeRF

Train a neural network to input position p
and direction d , output $\sigma(p)$ and $c(p, d)$

Training loss: Estimated pixel colors $C(r)$
should match actual pixel colors from images

NeRF: Network Architecture

Fully-connected network: Input position x and direction d , and output volume density and RGB color



Rather than pass raw xyz values to network,
instead use **positional encodings**:

$$\gamma(p) = (\sin(2^0 \pi p), \cos(2^0 \pi p), \dots, \sin(2^{L-1} \pi p), \cos(2^{L-1} \pi p))$$

NeRF: Network Architecture

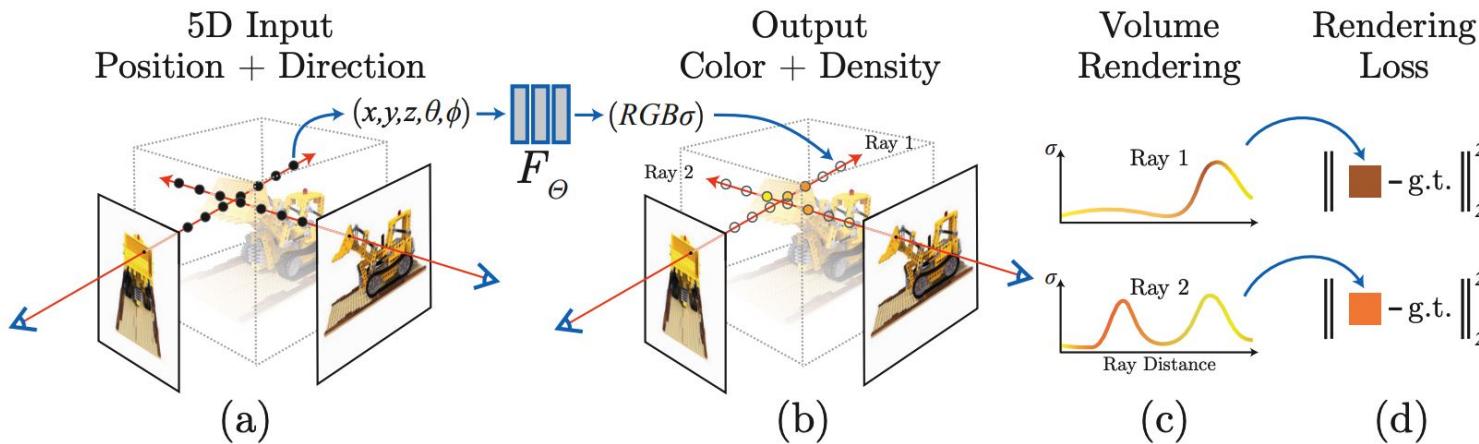


Fig. 2: An overview of our neural radiance field scene representation and differentiable rendering procedure. We synthesize images by sampling 5D coordinates (location and viewing direction) along camera rays (a), feeding those locations into an MLP to produce a color and volume density (b), and using volume rendering techniques to composite these values into an image (c). This rendering function is differentiable, so we can optimize our scene representation by minimizing the residual between synthesized and ground truth observed images (d).

NeRF



NeRF (2020)

Main Problem: Very slow!

Training: 1-2 days on a V100 GPU, for just a single scene!

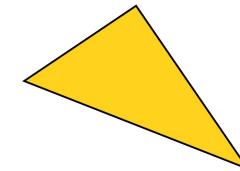
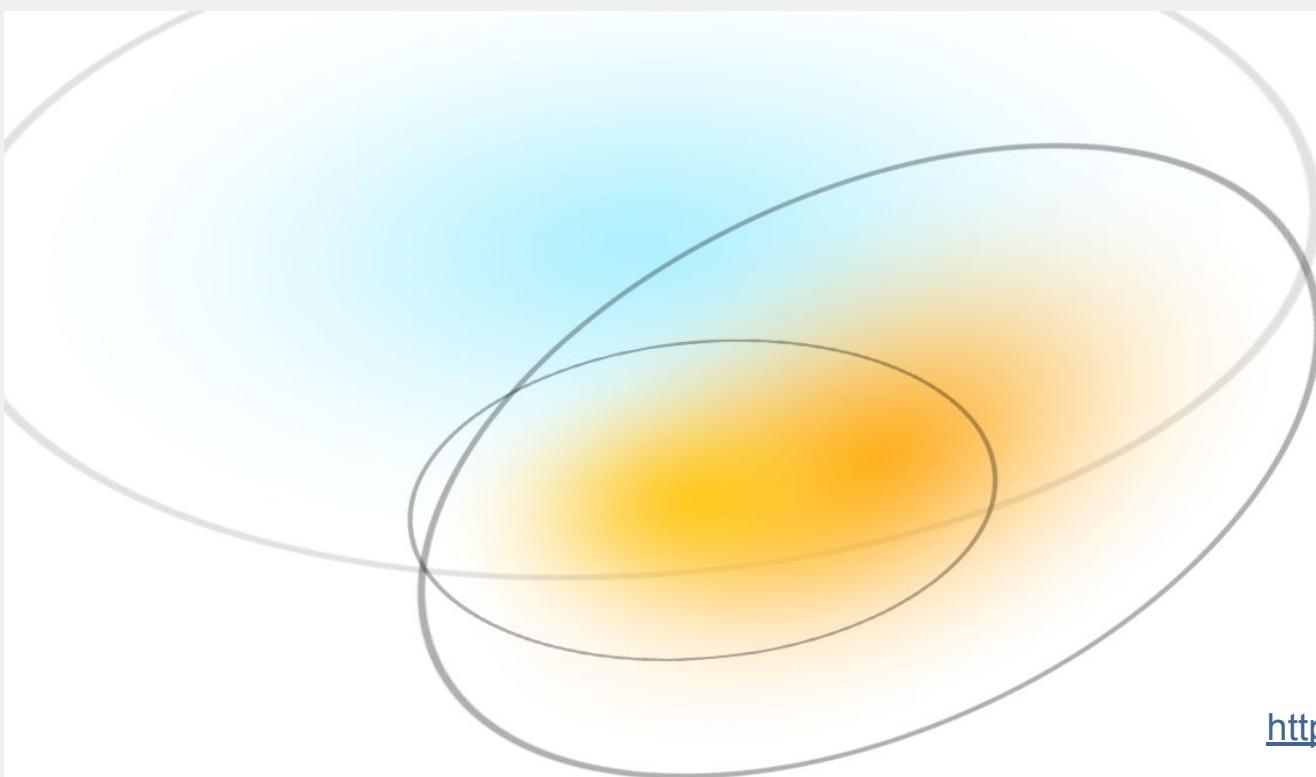
Inference: Sampling an image from a trained model:
 $(256 \times 256 \text{ pixels}) \times (224 \text{ samples per pixel})$
= 14.6M forward passes through MLP

Tons of follow-up work!

Cited by 11521

Gaussian Splatting (2023)

Recall: Triangle Meshes



[https://arxiv.org/pdf/2308.04079](https://arxiv.org/pdf/2308.04079.pdf)

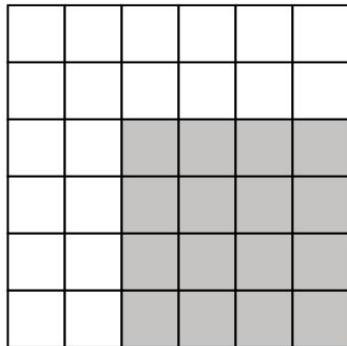
Gaussian Splatting



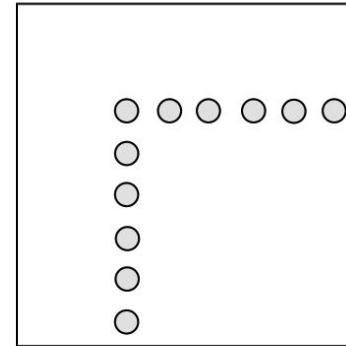
(More on this)

Summary

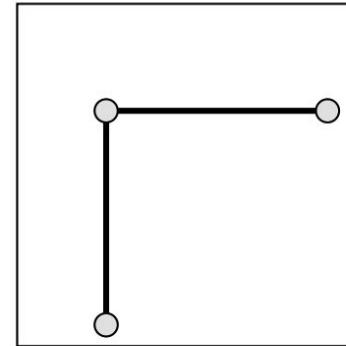
∞
∞
2
2
2
2
2



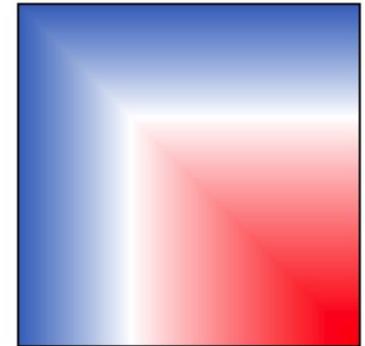
Depth
Map



Pointcloud



Mesh



Implicit
Surface

Reminder:

- **Final Project** “lightning talk” April 1, 2025 (tomorrow during discussion)
- **Canvas Quiz**

MCity Data Collection

Tomorrow (April 1st), 8am EST

Reach out to Cale Colony (GSI) ccolony@umich.edu

Dress warmly (26-44F)

<https://piazza.com/class/m4pgejar4ua2qf/post/228>

