

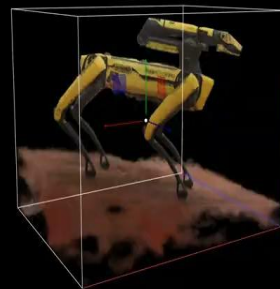
# DeepRob

[Student] Lecture 17

by *Michael Andrev, Chen Hu, Sean Coffey*

Implicit Surfaces, Geometry, NeRF

University of Michigan and University of Minnesota

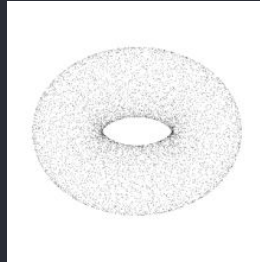




**Let's go over the basics!**

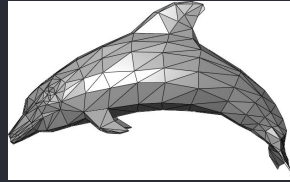
# Data Structures of 3D Geometry

## 1. Pointclouds



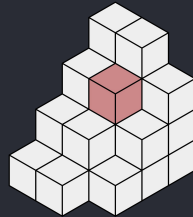
[https://commons.wikimedia.org/wiki/File:Point\\_cloud\\_torus.gif](https://commons.wikimedia.org/wiki/File:Point_cloud_torus.gif)

## 1. Meshes



[https://en.wikipedia.org/wiki/Polygon\\_mesh](https://en.wikipedia.org/wiki/Polygon_mesh)

## 1. Voxels



<https://en.wikipedia.org/wiki/Voxel>

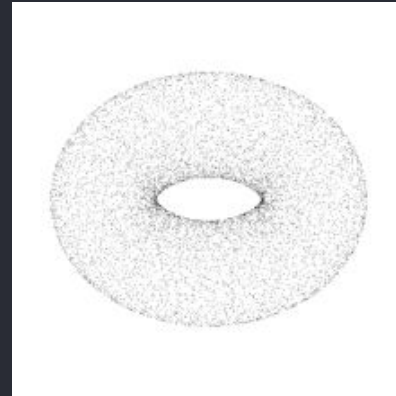
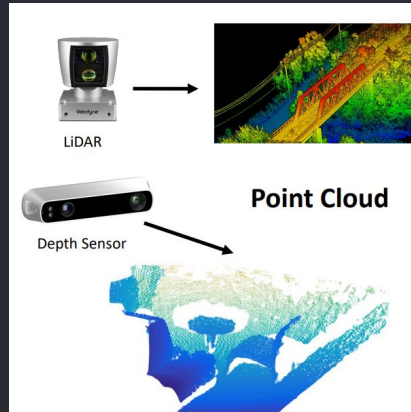
# Pointclouds

Point + Cloud = Pointcloud

A measurement unit that is represented using x, y, and z coordinates.

An aggregation of many small units of something.

A set of points in a space that represent some 3D shape or object



<https://geo-matching.com/content/how-to-get-the-best-precision-and-improve-pointcloud-accuracy>

[https://commons.wikimedia.org/wiki/File:Point\\_cloud\\_torus.gif](https://commons.wikimedia.org/wiki/File:Point_cloud_torus.gif)

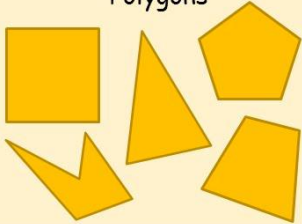
# Meshes

Meshes (or polygon meshes) are a 3D object that is composed of 1+ polygons

## What does polygon mean?

The word **polygon** comes from two Greek words.

### Polygons



The word **poly** means "many":  
**polymorph** = many forms  
**polymath** = someone with lots of different skills

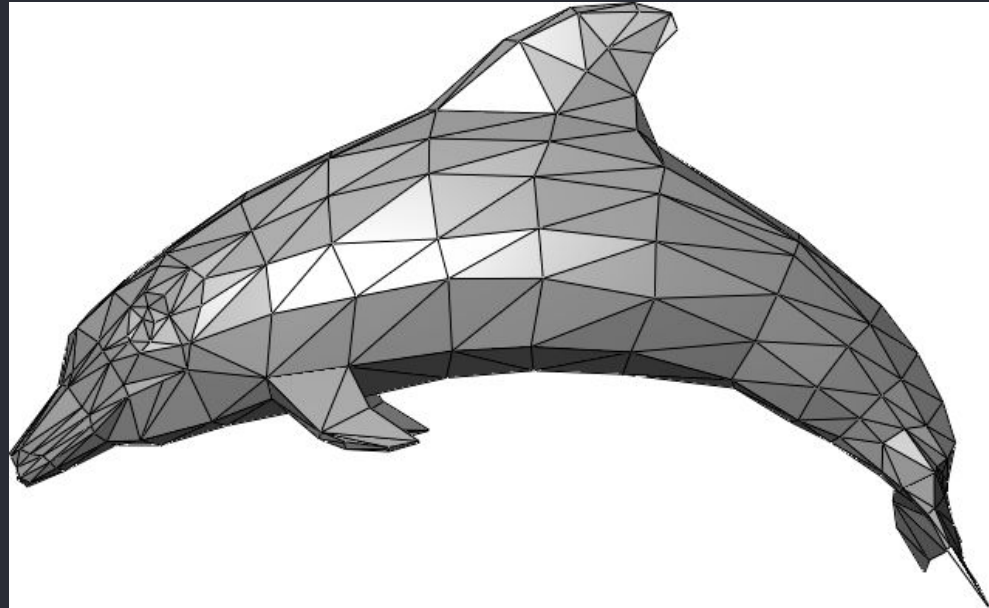
The word **gon** means "angle":  
**hexagon** = six angles

The word **polygon** means "many angles".

### Not polygons

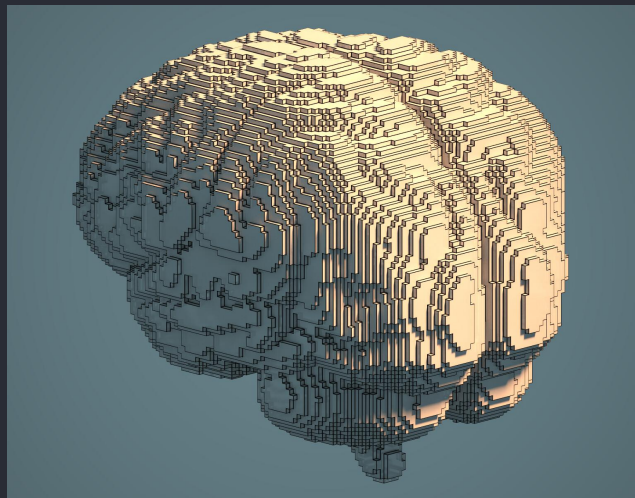


A more helpful definition:  
A polygon is a closed 2D shape made of straight lines (e.g. not a circle, shapes with curves or an open loop!).

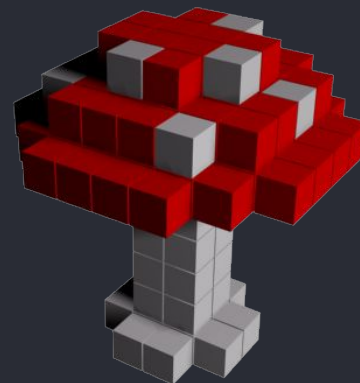


# Voxels

Voxels are 3D pixels represented as cubes.

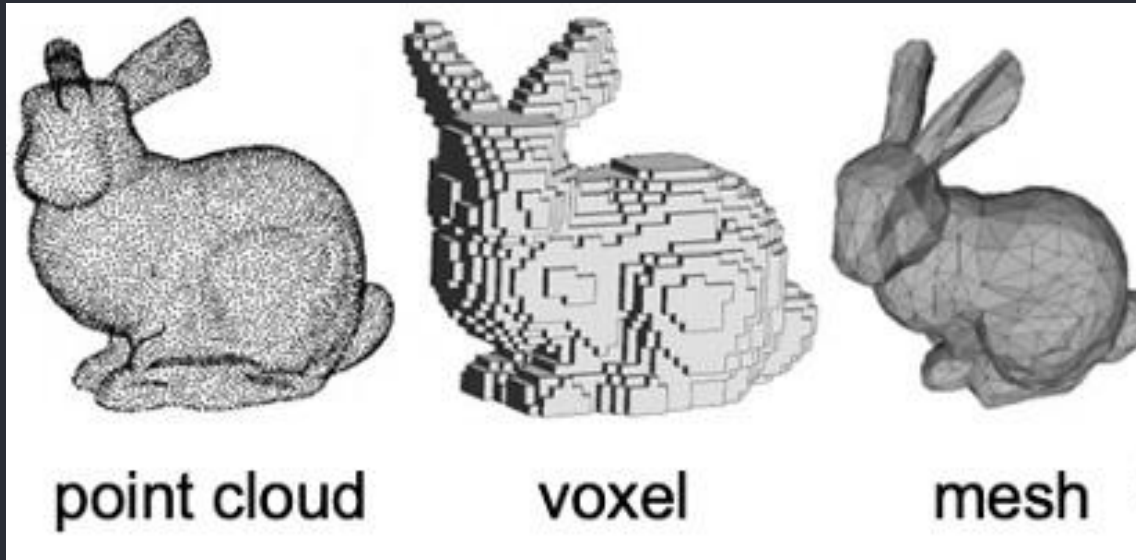


[https://www.google.com/url?sa=i&url=https%3A%2F%2Fdon\\_backos.artstation.com%2Fprojects%2FrR5Y8J&psig=AOvVaw1miitIbCKz2tjcvK0X7XW\\_&ust=1680288200369000&source=images&cd=vfe&ved=0CA8QjRxqFwoTCNicib\\_ChP4CFQAAAAAdAAAAABAD](https://www.google.com/url?sa=i&url=https%3A%2F%2Fdon_backos.artstation.com%2Fprojects%2FrR5Y8J&psig=AOvVaw1miitIbCKz2tjcvK0X7XW_&ust=1680288200369000&source=images&cd=vfe&ved=0CA8QjRxqFwoTCNicib_ChP4CFQAAAAAdAAAAABAD)



<https://blog.spatial.com/the-main-benefits-and-disadvantages-of-voxel-modeling>

# Summary of Data Structures of 3D Geometry





Questions?

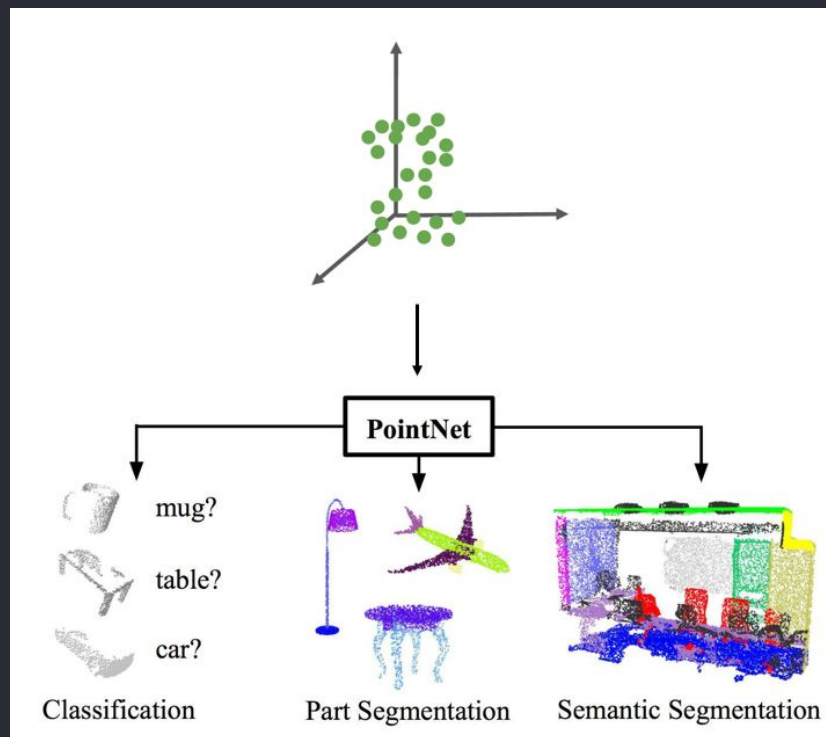


The background is a dark blue color with a network of thin, gold-colored lines forming various geometric shapes, primarily triangles and polygons, scattered across the frame. A prominent gold rectangular border is centered on the page, enclosing the main text.

# PointNet and PointNet++

# PointNet

1. **End-to-end learning** for irregular point data
2. **Unified** framework for various tasks



Qi, Charles R., et al. "Pointnet: Deep learning on point sets for 3d classification and segmentation." Proceedings of the IEEE conference on computer vision and pattern recognition. 2017.

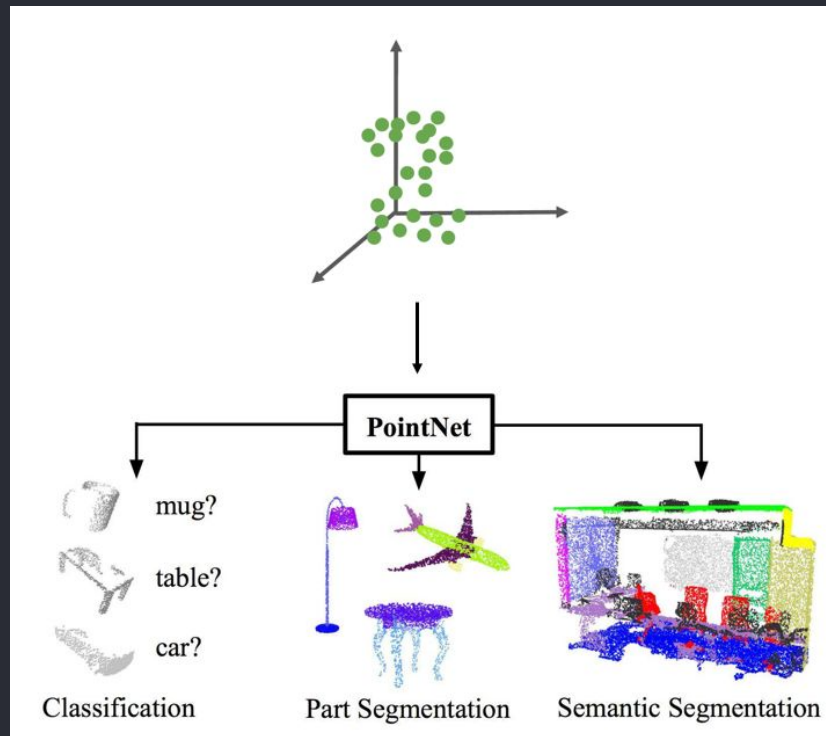
# PointNet

1. **End-to-end learning** for irregular point data
2. **Unified** framework for various tasks

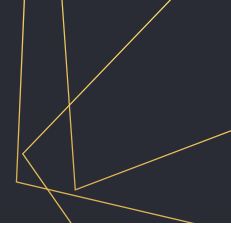
**PointNet** has to respect key characteristics for points clouds:

1. Point Permutation Invariance
2. Spatial Transformation Invariance
3. Sampling Invariance

Qi, Charles R., et al. "Pointnet: Deep learning on point sets for 3d classification and segmentation." Proceedings of the IEEE conference on computer vision and pattern recognition. 2017.



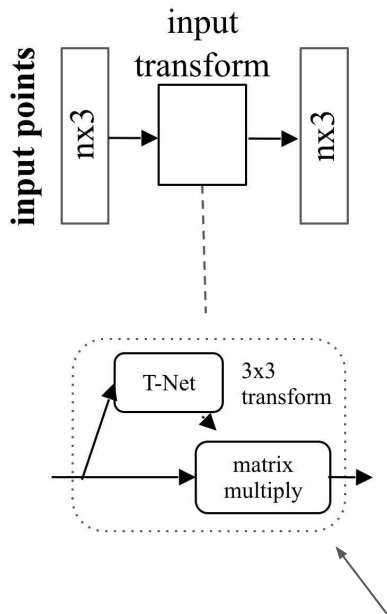
# PointNet Classification Network



input points

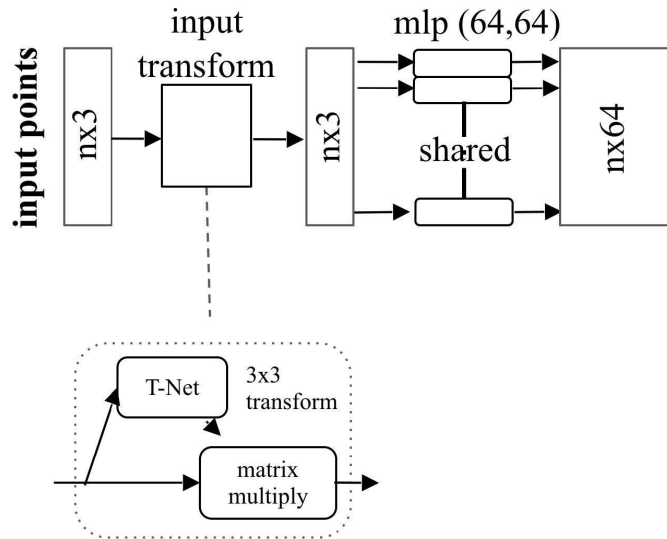
$n \times 3$

# PointNet Classification Network

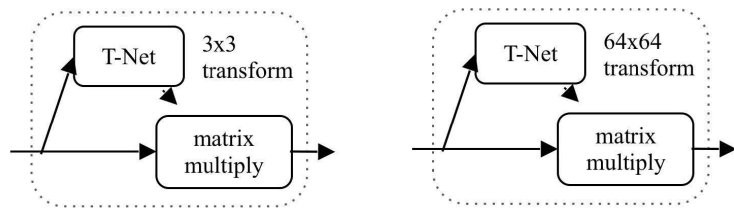
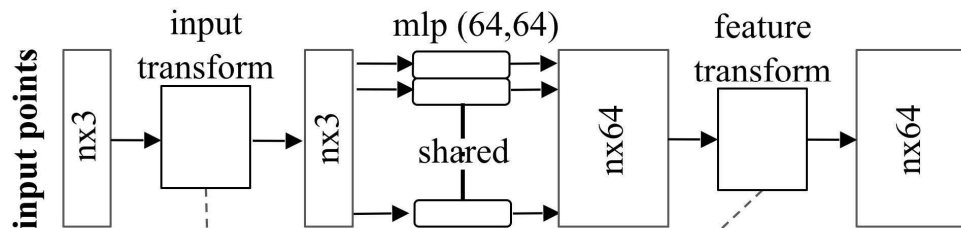


Encode Spatial Transformation Invariance

# PointNet Classification Network



# PointNet Classification Network

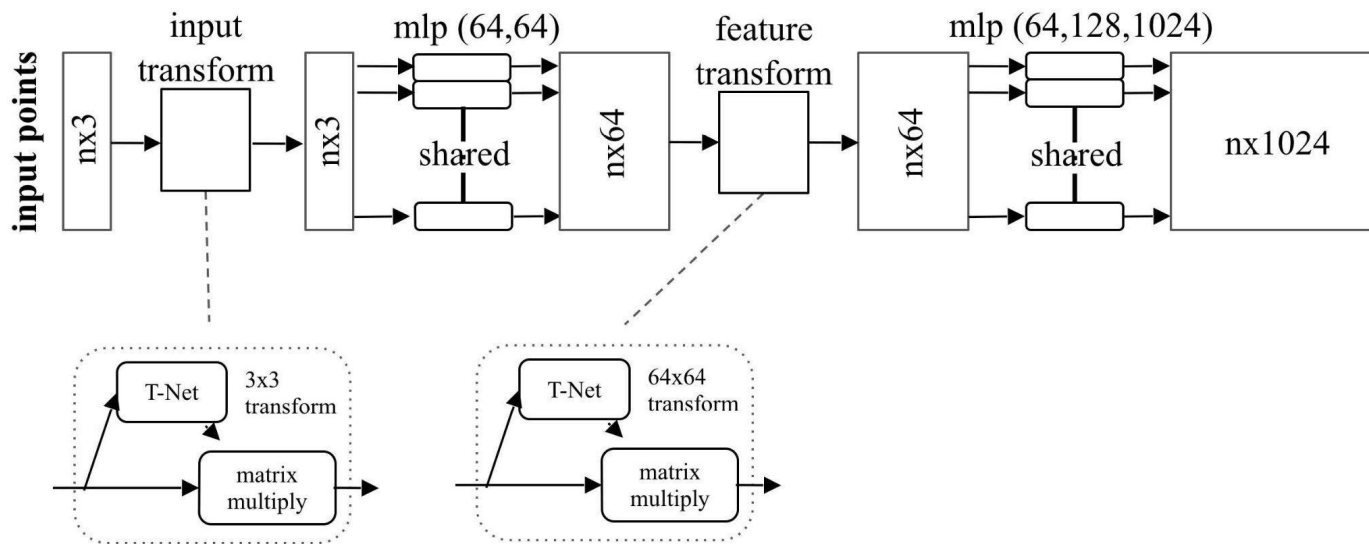


Encode Spatial Transformation Invariance

**Regularization loss:**

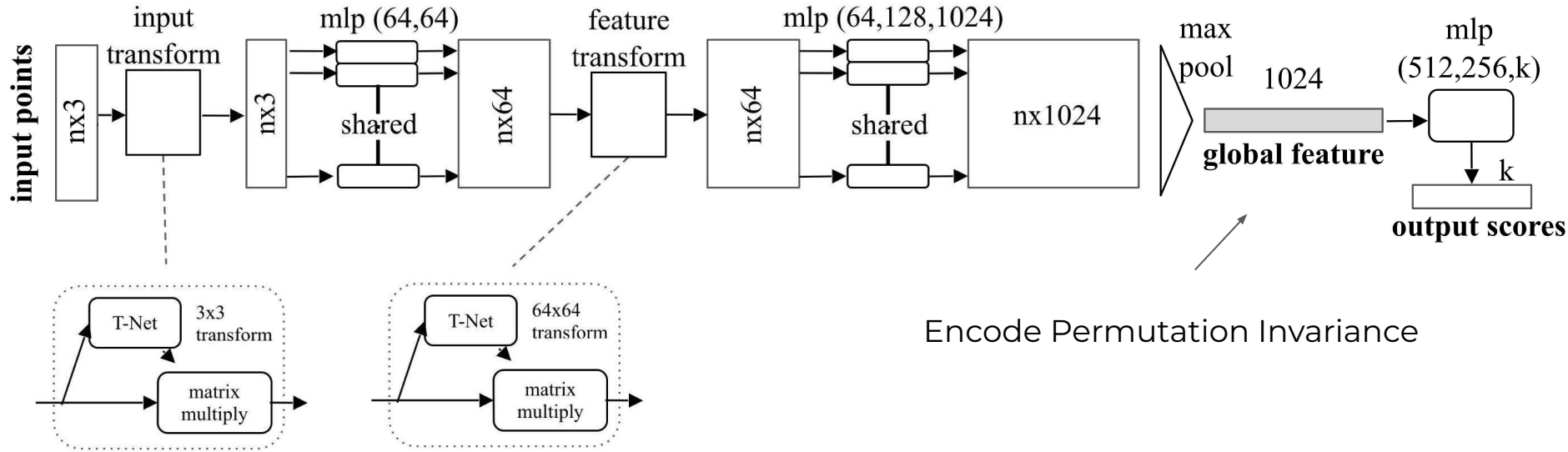
Transform matrix close to orthogonal:  $L_{reg} = \|I - AA^T\|_F^2$

# PointNet Classification Network

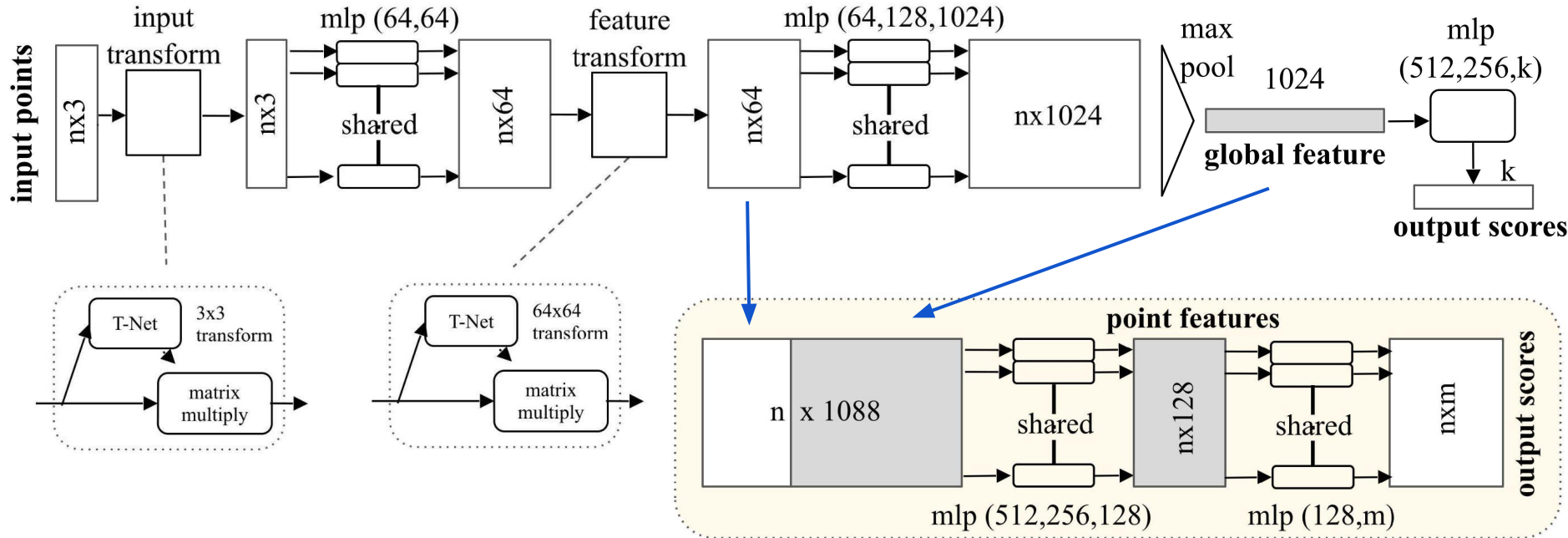




# PointNet Classification Network



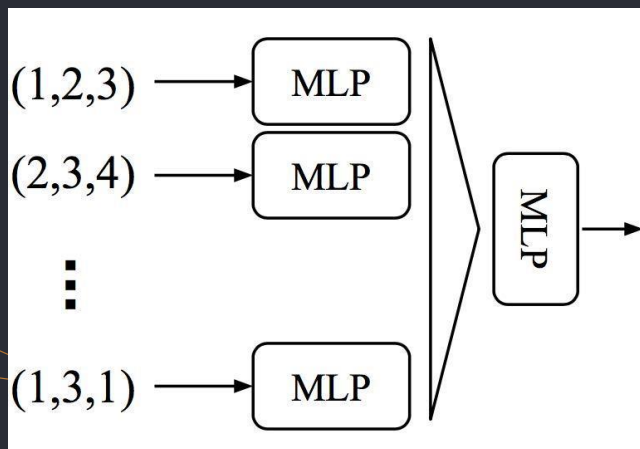
# PointNet Segmentation Network



# PointNet++

Limitation of PointNet - Global feature learning

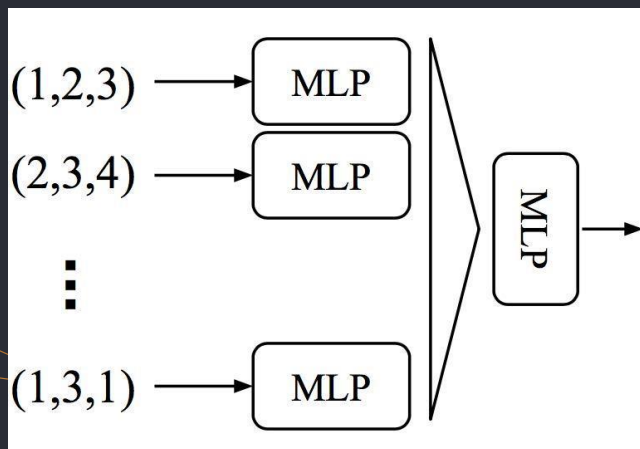
1. No local context
2. Limited local invariance



# PointNet++

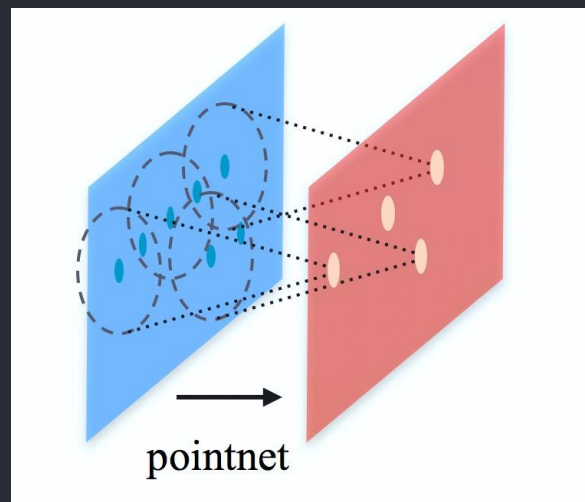
Limitation of PointNet - Global feature learning

1. No local context
2. Limited local invariance

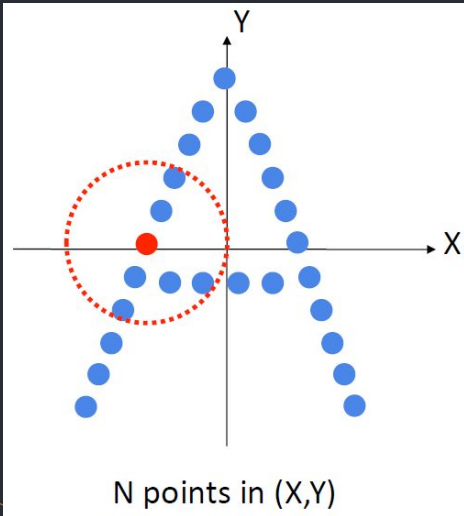


PointNet++: recursively apply pointnet at local regions:

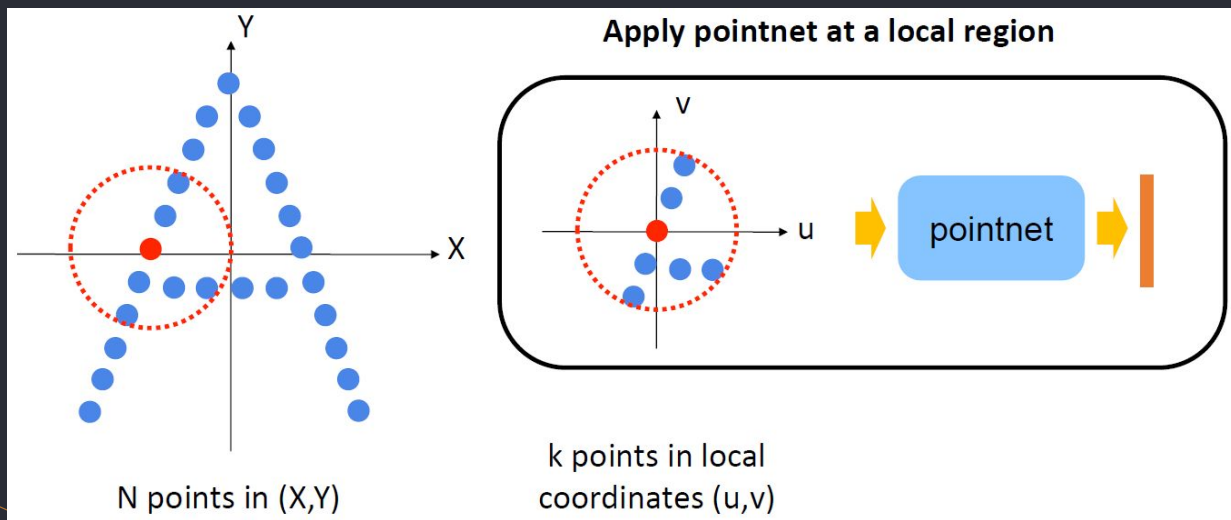
1. Hierarchical feature learning
2. Local translational invariance
3. Permutation invariance



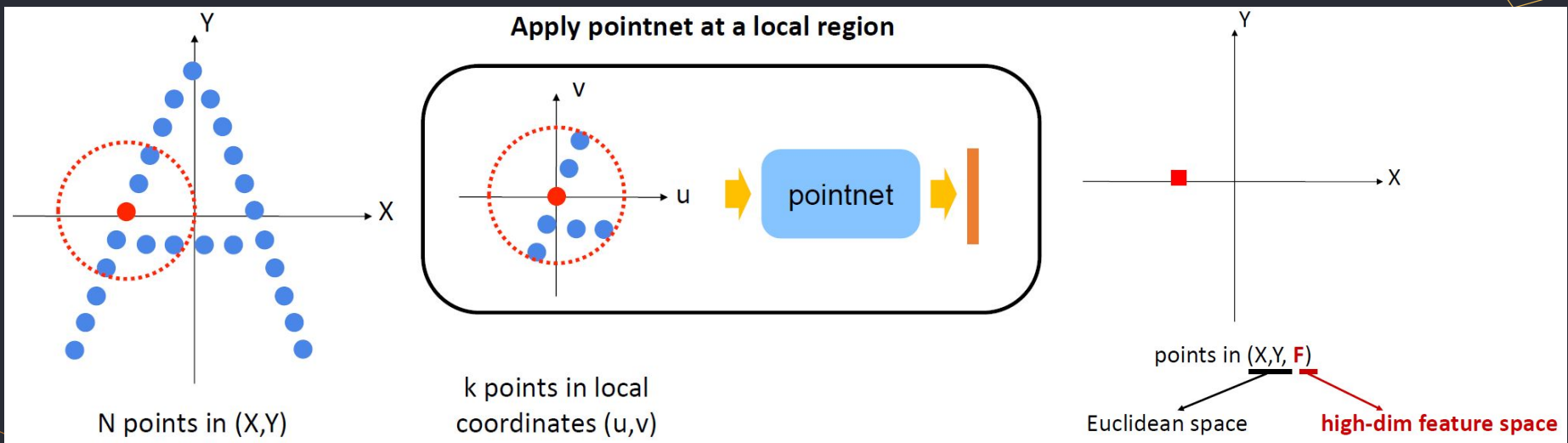
# Hierarchical Point Feature Learning



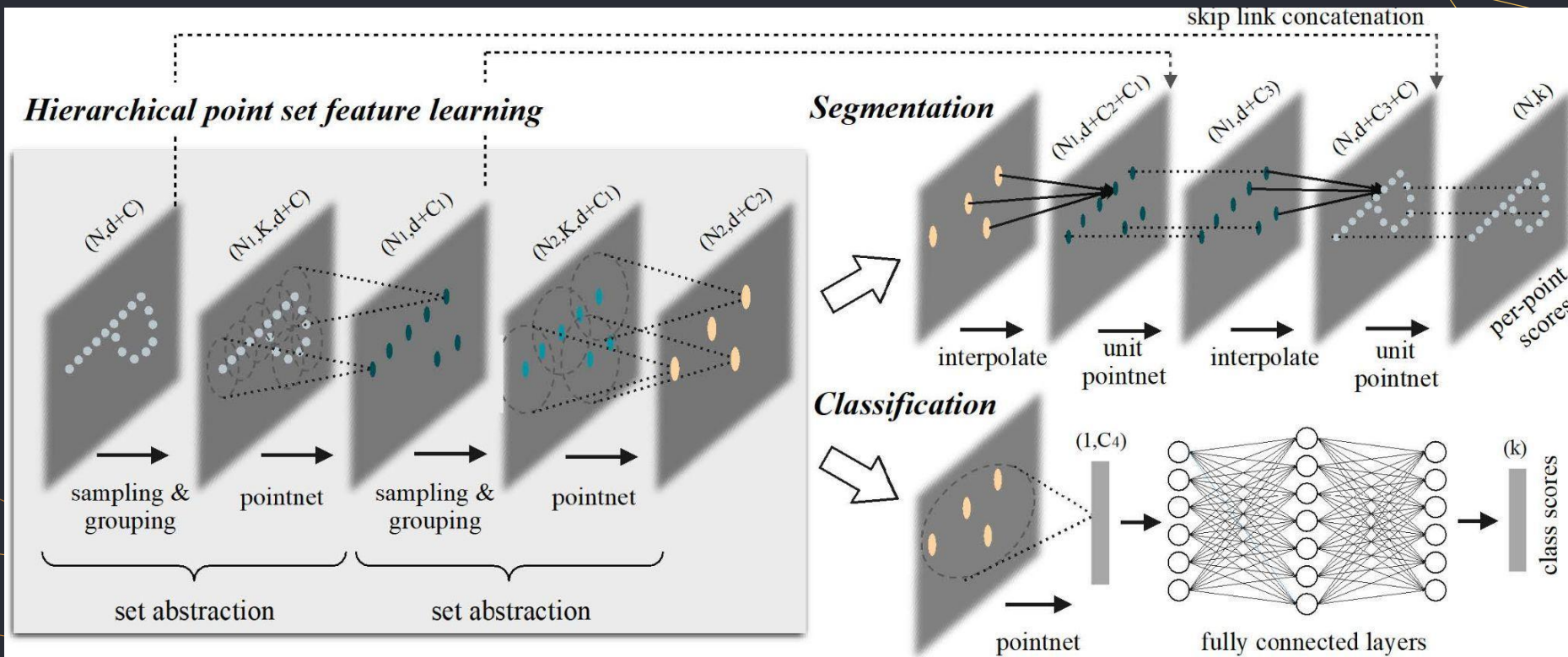
# Hierarchical Point Feature Learning



# Hierarchical Point Feature Learning

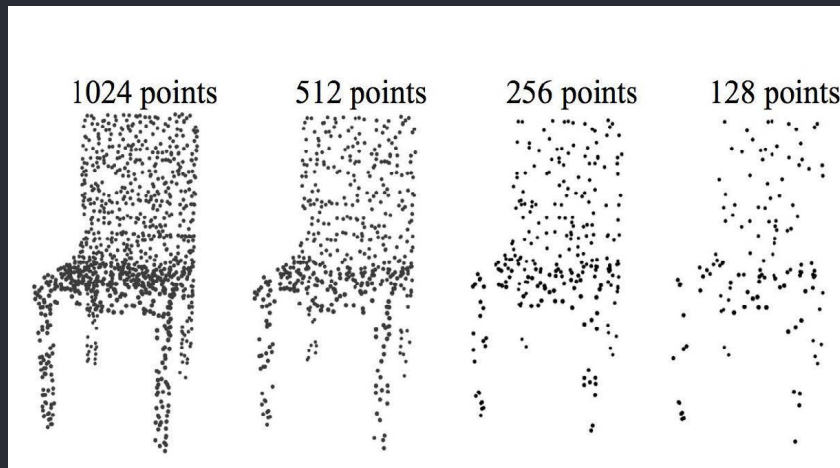


# PointNet++

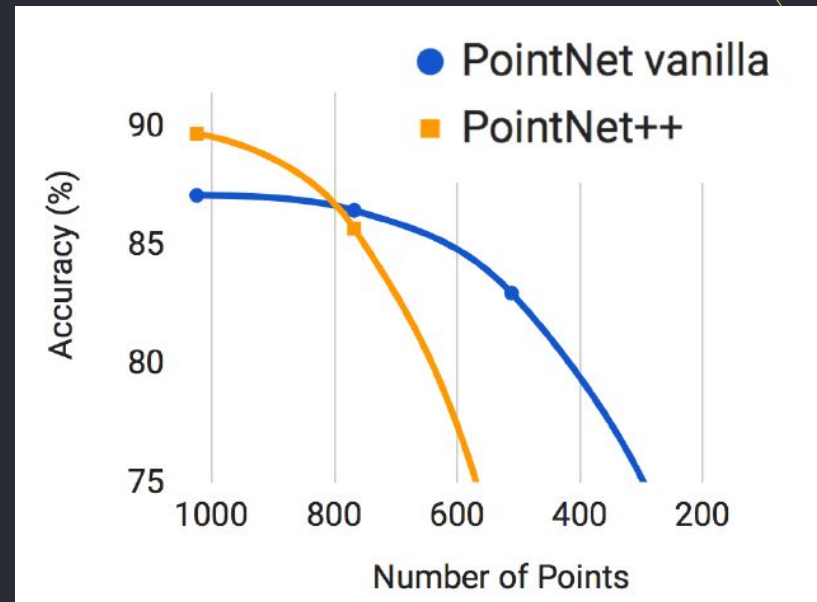




# PointNet and PointNet++ Comparison

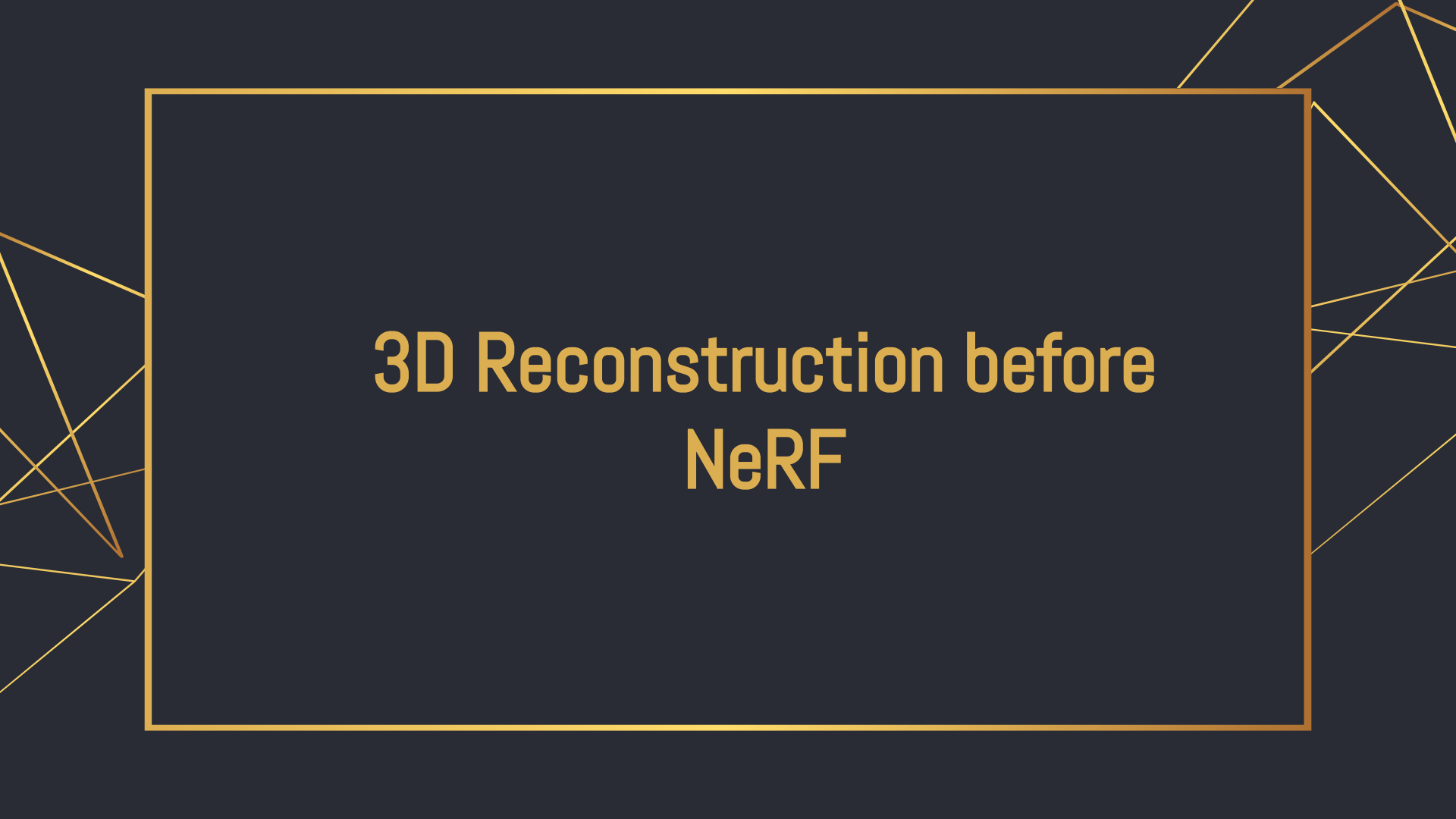


Density Variation





Questions?



# 3D Reconstruction before NeRF

# Reconstruction of 3D scenes and Geometries

The process of capturing the shape and appearance of real objects<sup>1</sup>

## Active Methods

- Actively interfere with the reconstructed object, either mechanically or radiometrically (ex laser range finder, 3D ultrasound)<sup>2</sup>

## Passive Methods

- Measure the radiance reflected or emitted by the object's surface to infer its 3D structure through image understanding<sup>3</sup>

<sup>1</sup>Moons, Theo, Luc Van Gool, and Maarten Vergauwen. "3D reconstruction from multiple images part 1: Principles."

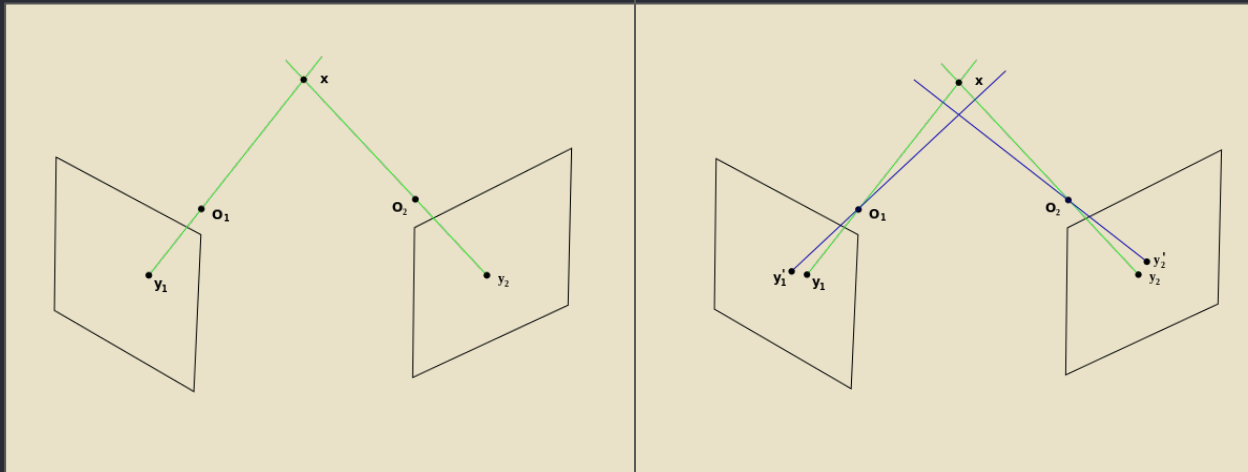
<sup>2</sup>Mahmoudzadeh, Ahmadreza; Golroo, Amir; Jahanshahi, Mohammad R.; Firoozi Yeganeh, Sayna (January 2019). "Estimating Pavement Roughness by Fusing Color and Depth Data Obtained from an Inexpensive RGB-D Sensor"

<sup>3</sup>Buelthoff, Heinrich H., and Alan L. Yuille. "Shape-from-X: Psychophysics and computation Archived 2011-01-07 at the Wayback Machine."

# Triangulation

The process of determining a point in 3D space given its projections onto two, or more, images.

In order to solve this problem it is necessary to know the parameters of the camera projection function from 3D to 2D for the cameras involved, in the simplest case represented by the camera matrices.<sup>1</sup>

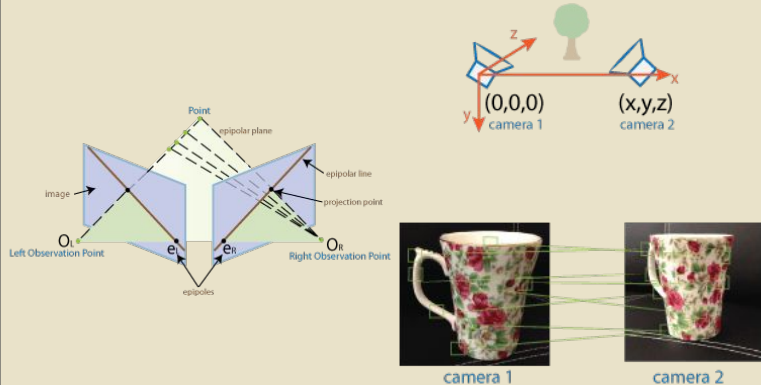


<sup>1</sup>Richard Hartley and Andrew Zisserman (2003). "Multiple View Geometry in computer vision"

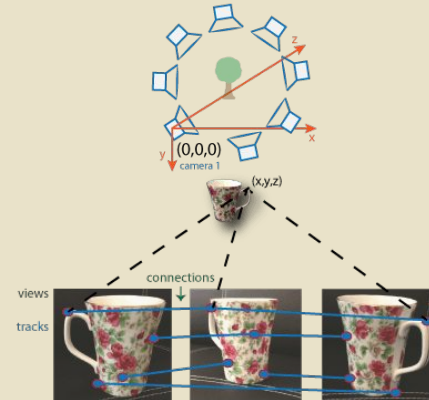
# Structure from Motion (SfM)

Structure from motion (SfM) is the process of estimating the 3-D structure of a scene from a set of 2-D images. SfM is used in many applications, such as 3-D scanning, augmented reality, and visual simultaneous localization and mapping (vSLAM)

## Structure from Motion from Two Views

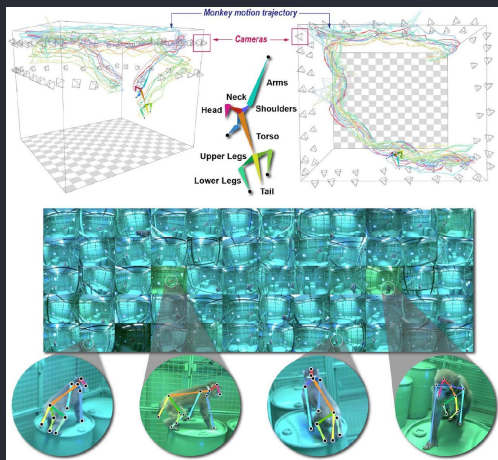


## Structure from Motion from Multiple Views

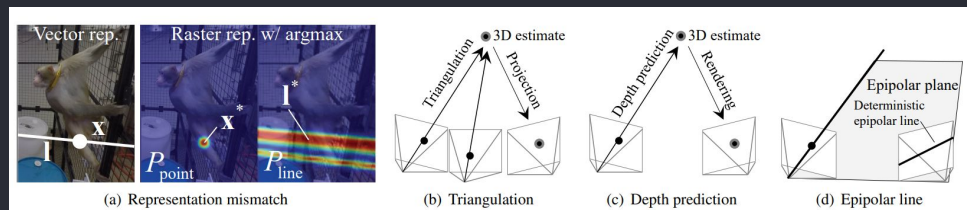


# Simultaneous Localization and Mapping (SLAM)

The computational problem of constructing or updating a map of an unknown environment while simultaneously keeping track of an agent's location within it.



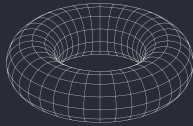
## Open Monkey Studio



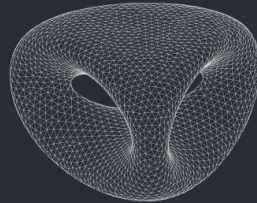
# Implicit Surfaces

An implicit surface is the set of zeros of a function of three variables. Implicit means that the equation is not solved for  $x$  or  $y$  or  $z$ .

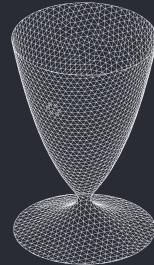
$$F(x, y, z) = 0.$$



$$(x^2 + y^2 + z^2 + R^2 - a^2)^2 - 4R^2(x^2 + y^2) = 0.$$



$$2y(y^2 - 3x^2)(1 - z^2) + (x^2 + y^2)^2 - (9z^2 - 1)(1 - z^2) = 0$$

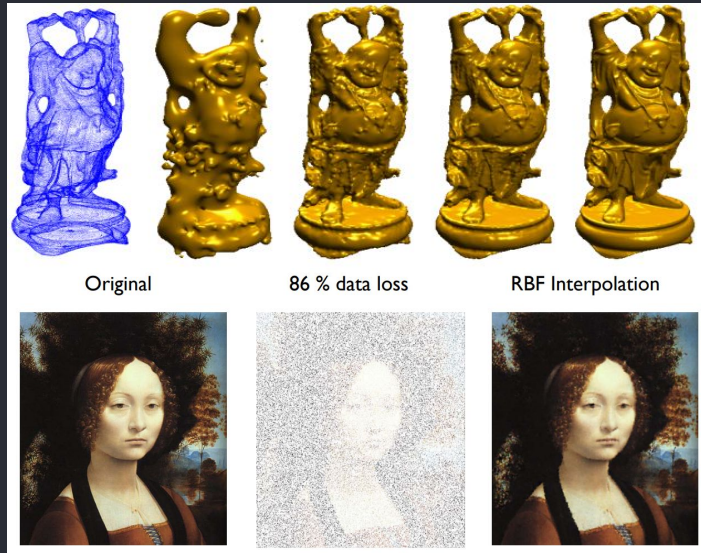


$$x^2 + y^2 - (\ln(z + 3.2))^2 - 0.02 = 0$$

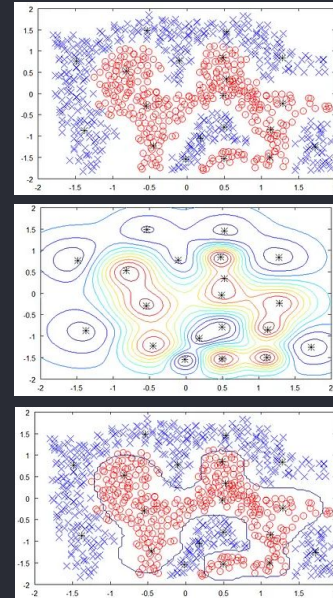


# Radial Basis Function

A radial basis function (RBF) is a real-valued function whose value depends only on the distance between the input and some fixed point

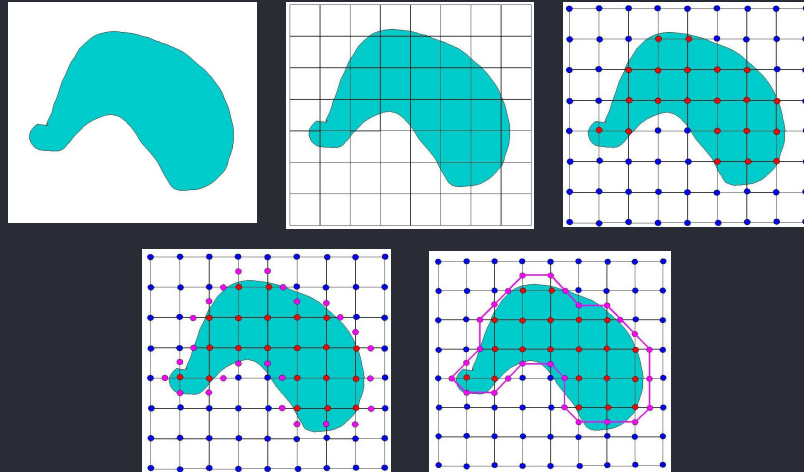


Images from: <https://www.cs.jhu.edu/~misha/Fall05/Papers/carr01.pdf>



# Marching Cubes Algorithm

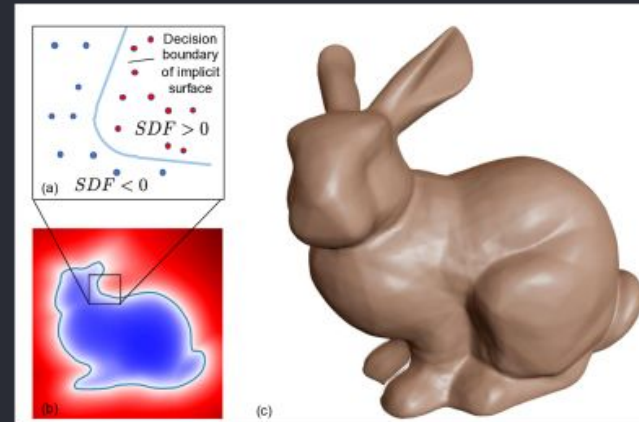
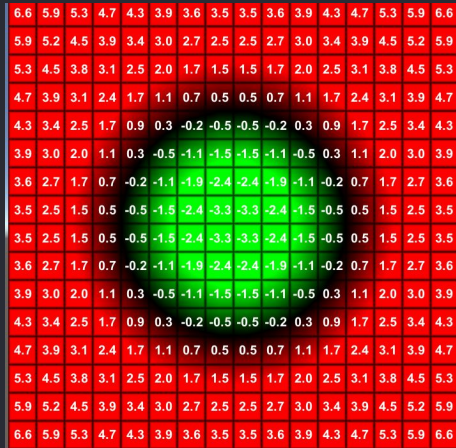
Marching cubes is a computer graphics algorithm, published in the 1987 SIGGRAPH proceedings by Lorensen and Cline, for extracting a polygonal mesh of an isosurface from a three-dimensional discrete scalar field (the elements of which are sometimes called voxels)



Images from:  
[https://www.cs.carleton.edu/cs\\_compos/0405/shape/marching\\_cubes.html](https://www.cs.carleton.edu/cs_compos/0405/shape/marching_cubes.html)

# Signed Distance Function

The signed distance function is a mathematical function that is used to define implicit surfaces.





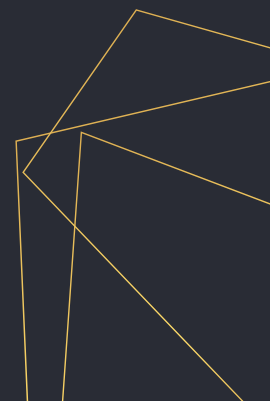
Questions?



**Finally, NeRF... Well... Almost!**

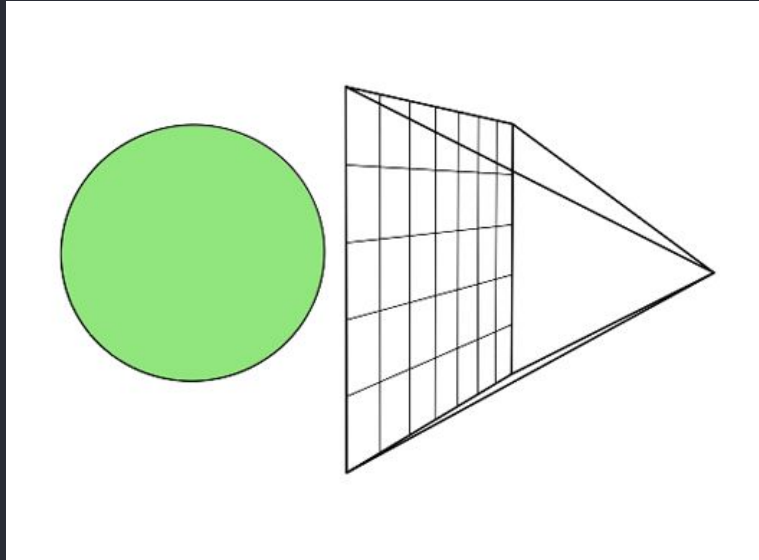


# Computer Graphic concepts

1. Rendering
  2. Volume Rendering
  3. Volume Synthesis
- 

# Rendering

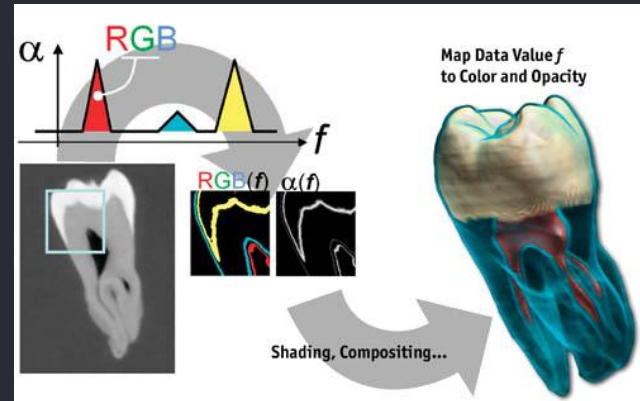
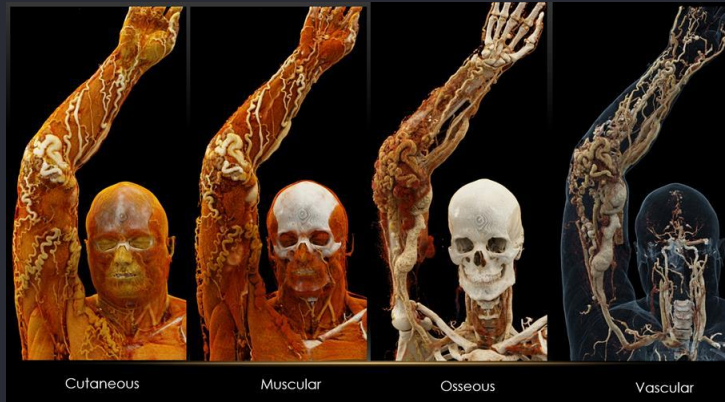
Generating an image (render) from a 2D/3D model.



# Volume Rendering

Create a 2D projection of 3D voxel data

3D objects  $\rightarrow$  2D images





# Volume Synthesis

Create a 3D view from a 2D scale

2D images  $\rightarrow$  3D objects



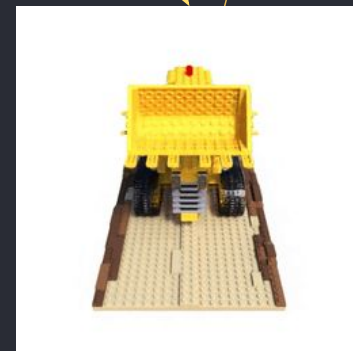
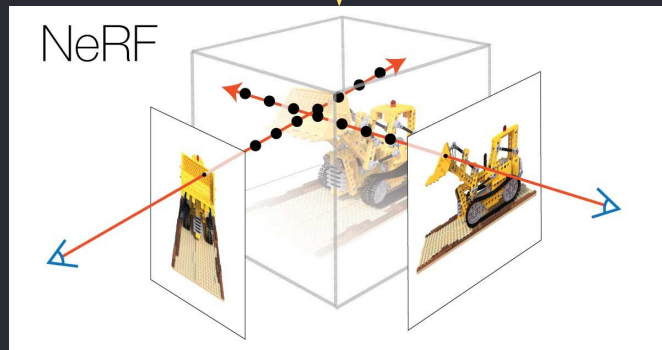
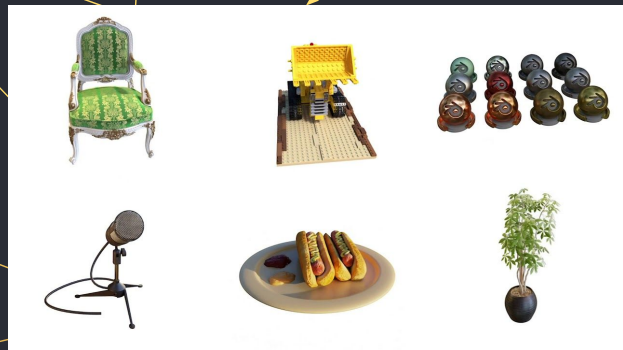


Questions?



**Finally, NeRF... For Real!**

# NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis

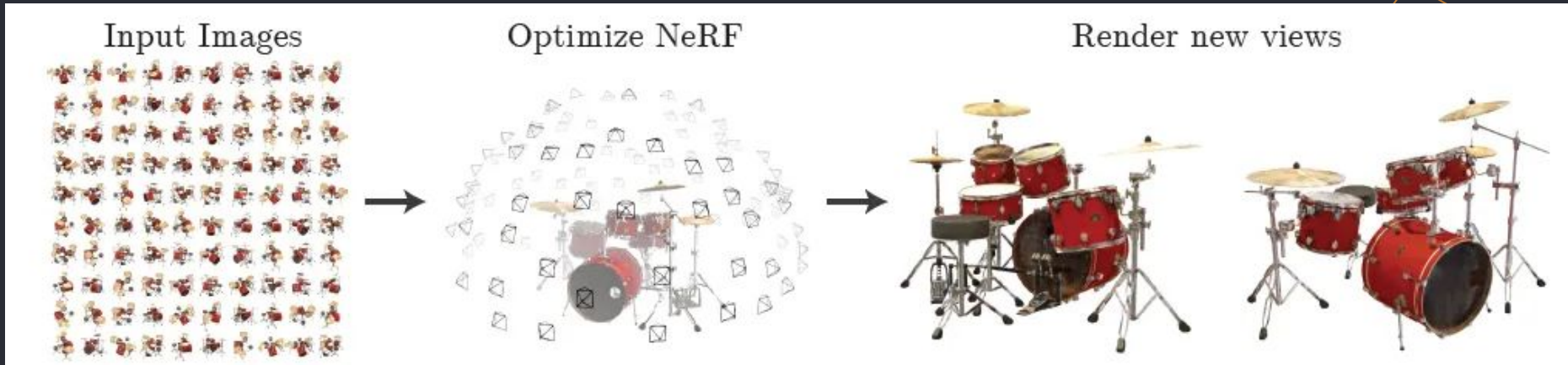


Paper Authors: Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, Ren Ng

# What is a NeRF?

Technical Definition: Fully-connected neural network that can generate new views of 3D scenes based on a set of 2D images.

**\* Think of NeRF as volume rendering + coordinate-based network!**



## Q. What are we trying to do?

- A. We are trying to generate a new picture/view of an object by existing sparse image inputs

## Q. How do we do that?

- A. Training a neural network which representing the 3D scene



# NeRF View-Dependent Appearance



# Geometry Visualization



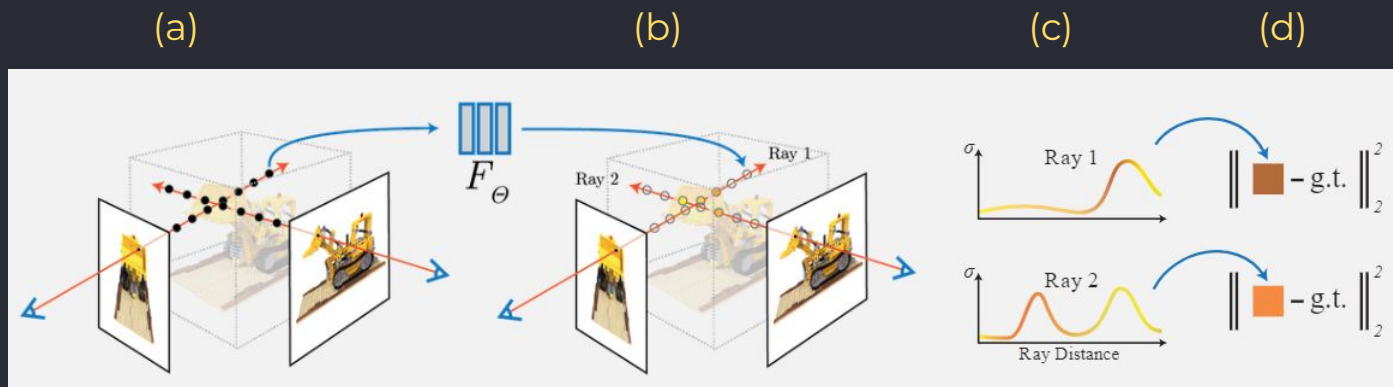
<https://www.matthewtancik.com/nerf>





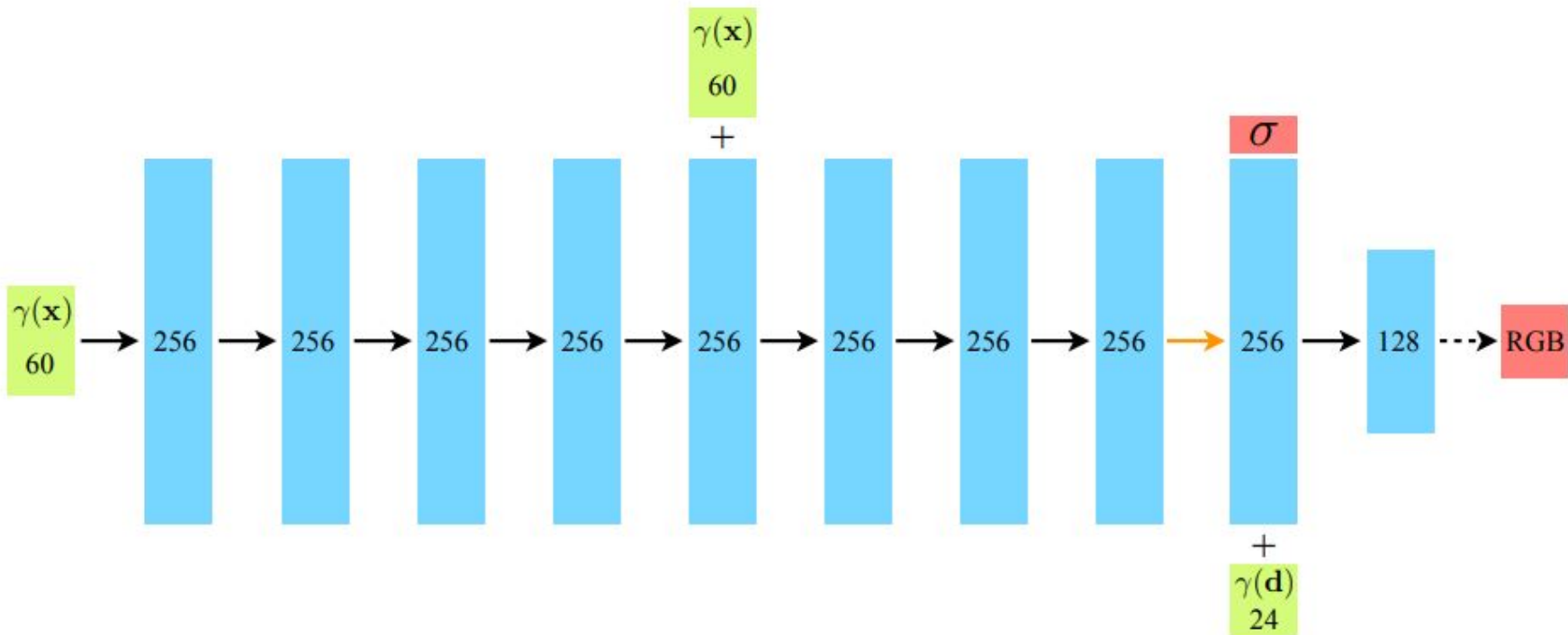
**Let's dive into the concepts  
surrounding NeRF!**

# Brief Overview of NeRF



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

**We overfit our MLP to this scene (this is VERY unusual in classical deep learning)!**



# Volume Rendering with Radiance Fields

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt, \text{ where } T(t) = \exp\left(-\int_{t_n}^t \sigma(\mathbf{r}(s)) ds\right)$$

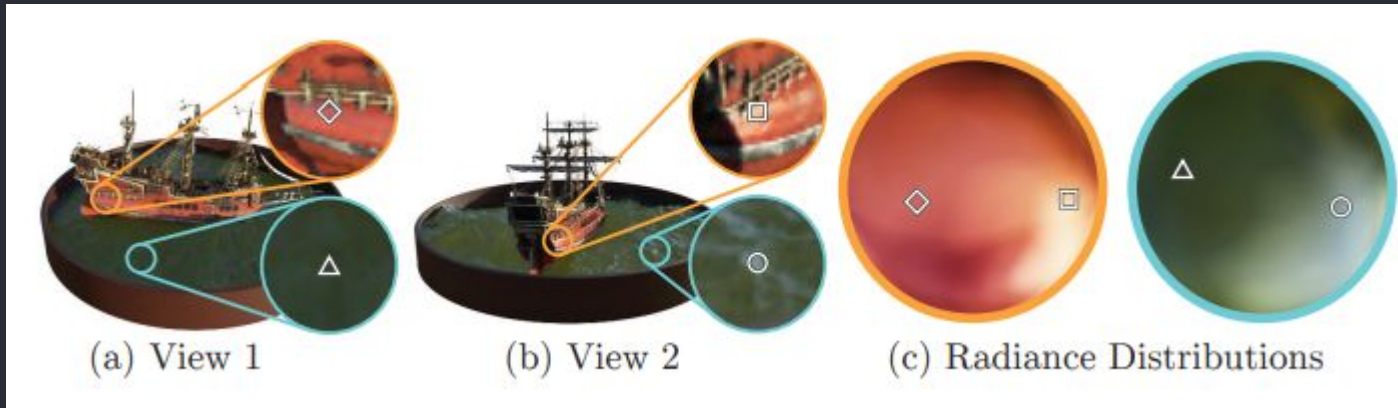
Expected  
color of  
camera ray

Density at  
the point  
(in terms of  
occlusion)

How much light has been  
blocked up to point  $t$ ?

Color at point  $\mathbf{r}(t)$   
from viewing  
direction  $\mathbf{d}$

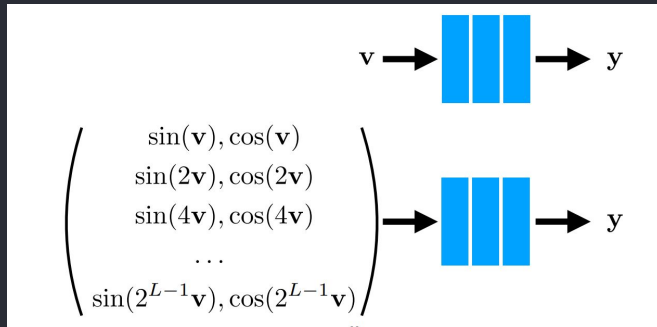
# View-Dependent Emitted Radiance



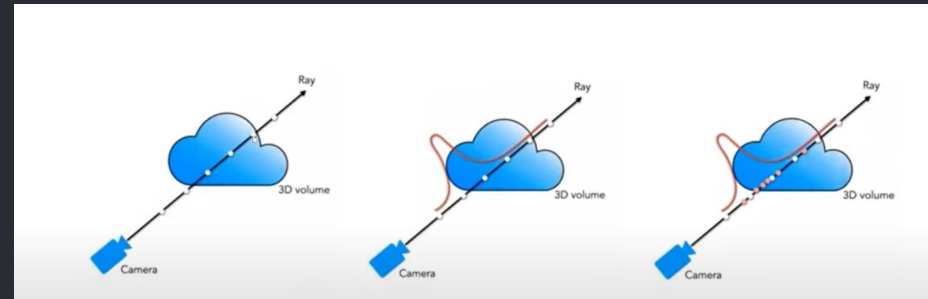
<https://arxiv.org/pdf/2003.08934.pdf>

# Optimizing a Neural Radiance Field

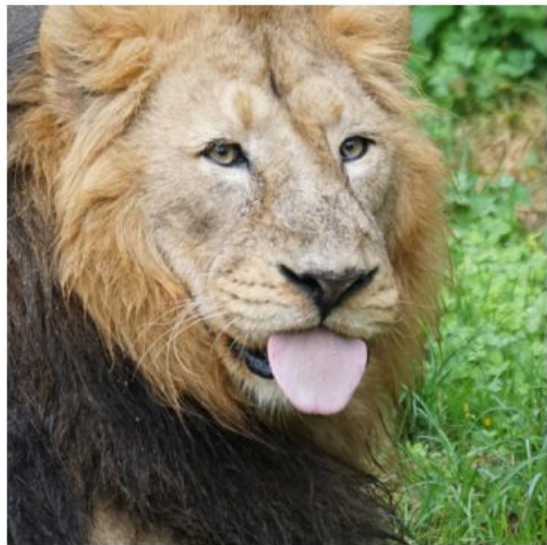
## 1. Positional Encoding



## 2. Hierarchical Volume Sampling



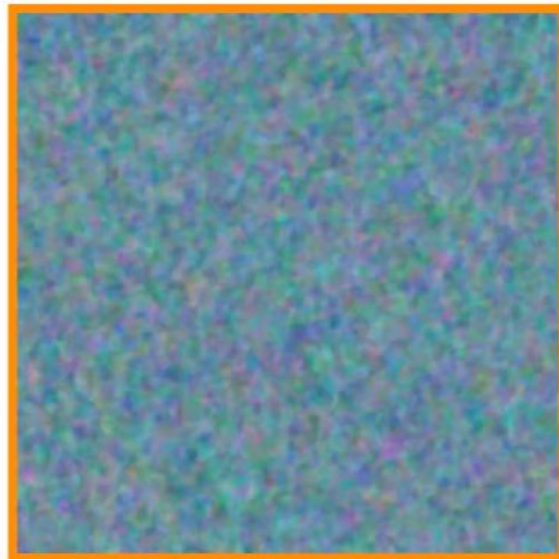
# Why do we use Positional Encoding?



Ground truth image



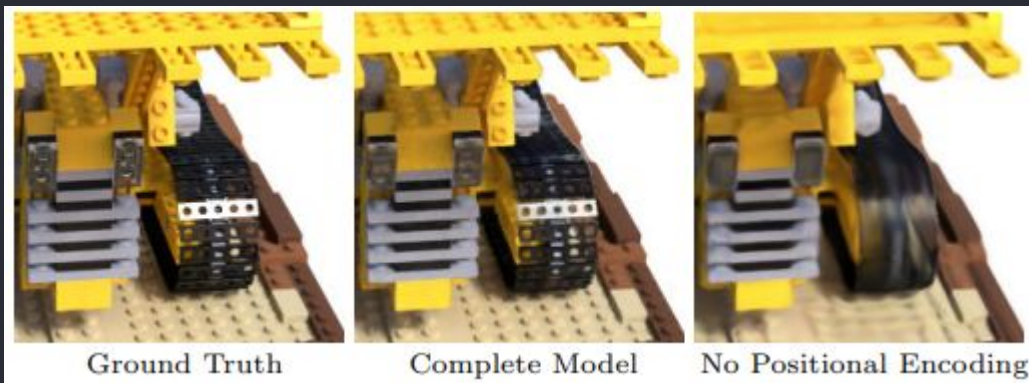
Neural network output **without**  
high frequency mapping



Neural network output **with**  
high frequency mapping



# Positional Encoding

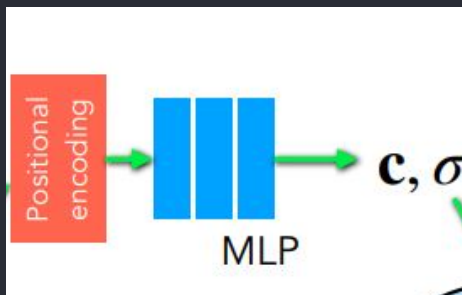


$$F_{\Theta} = F'_{\Theta} \circ \gamma$$

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

Compositions of a regular MLP and the mapping to a high-dimensional space

Mapping from  $\mathbb{R}$  to to  $\mathbb{R}^{2L}$

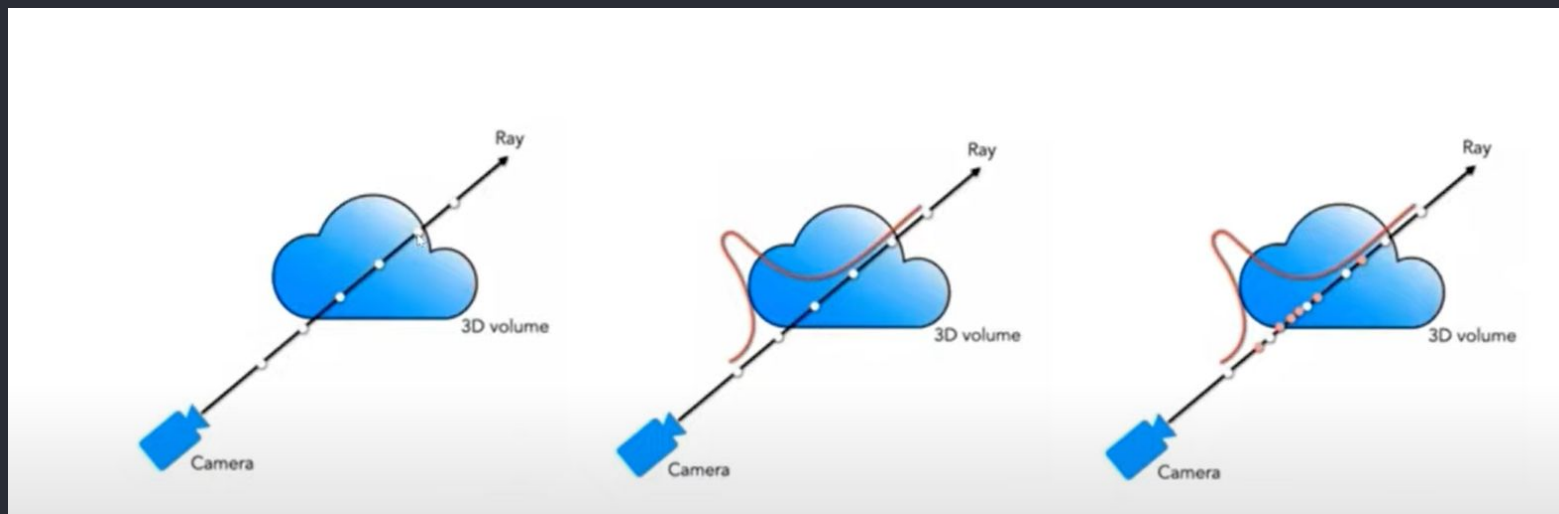


<https://arxiv.org/pdf/2003.08934.pdf>

<http://graphics.stanford.edu/courses/cs348n-22-winter/>

# Hierarchical Volume Sampling

$$\hat{C}_c(\mathbf{r}) = \sum_{i=1}^{N_c} w_i c_i, \quad w_i = T_i(1 - \exp(-\sigma_i \delta_i))$$

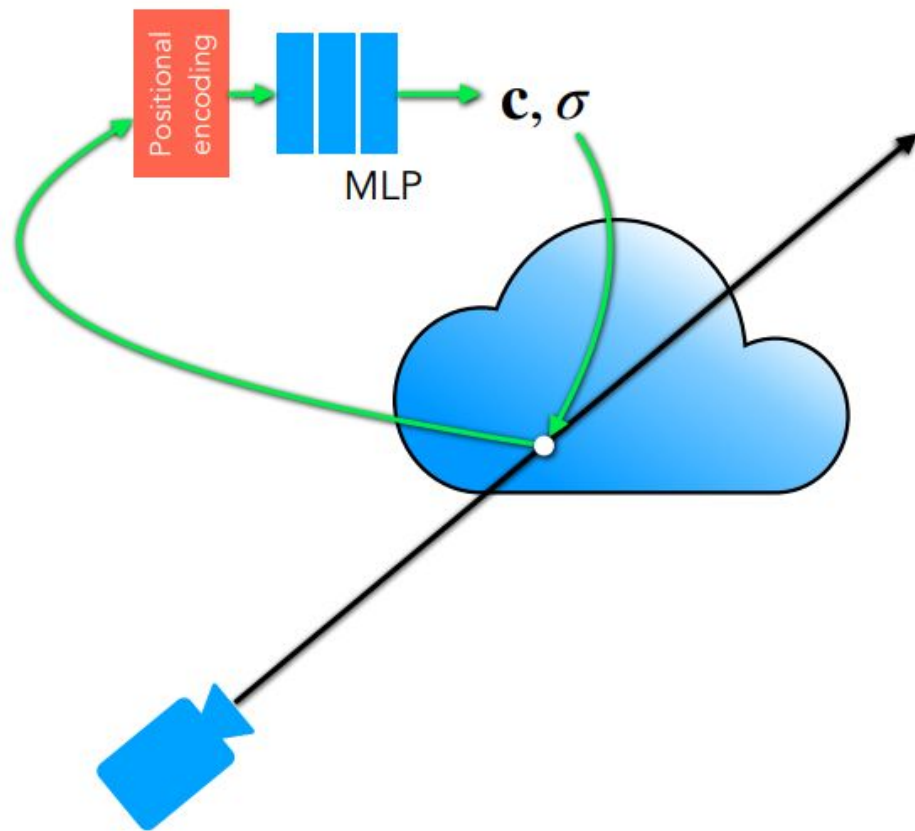


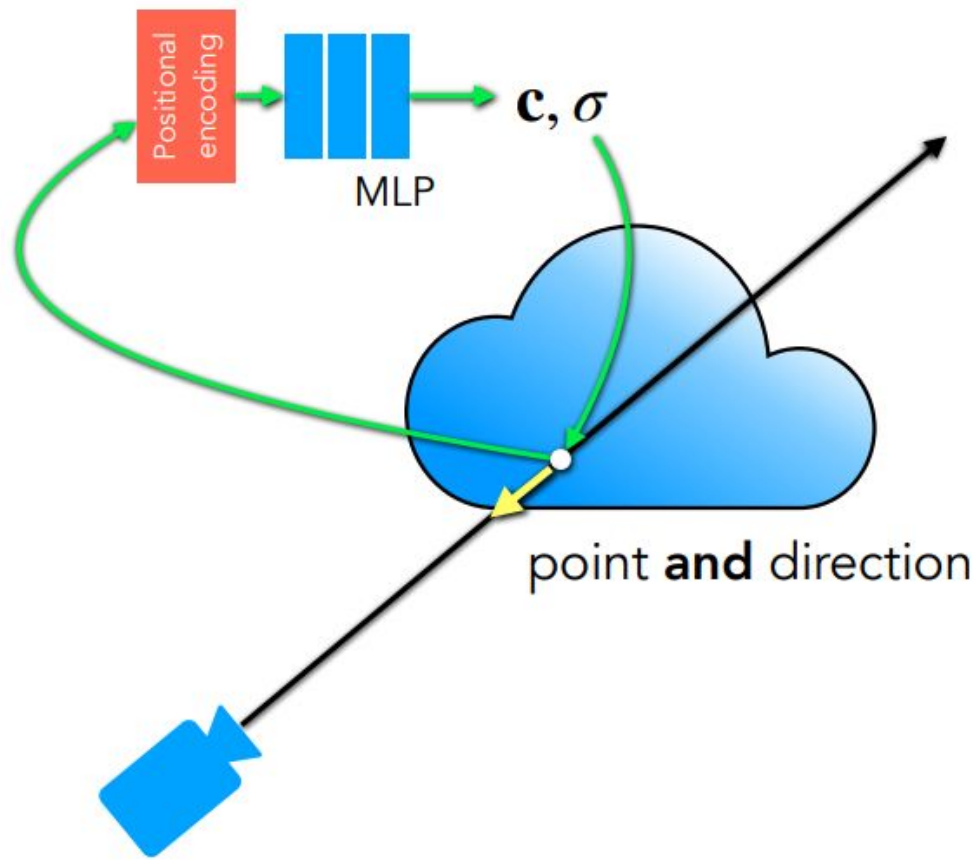


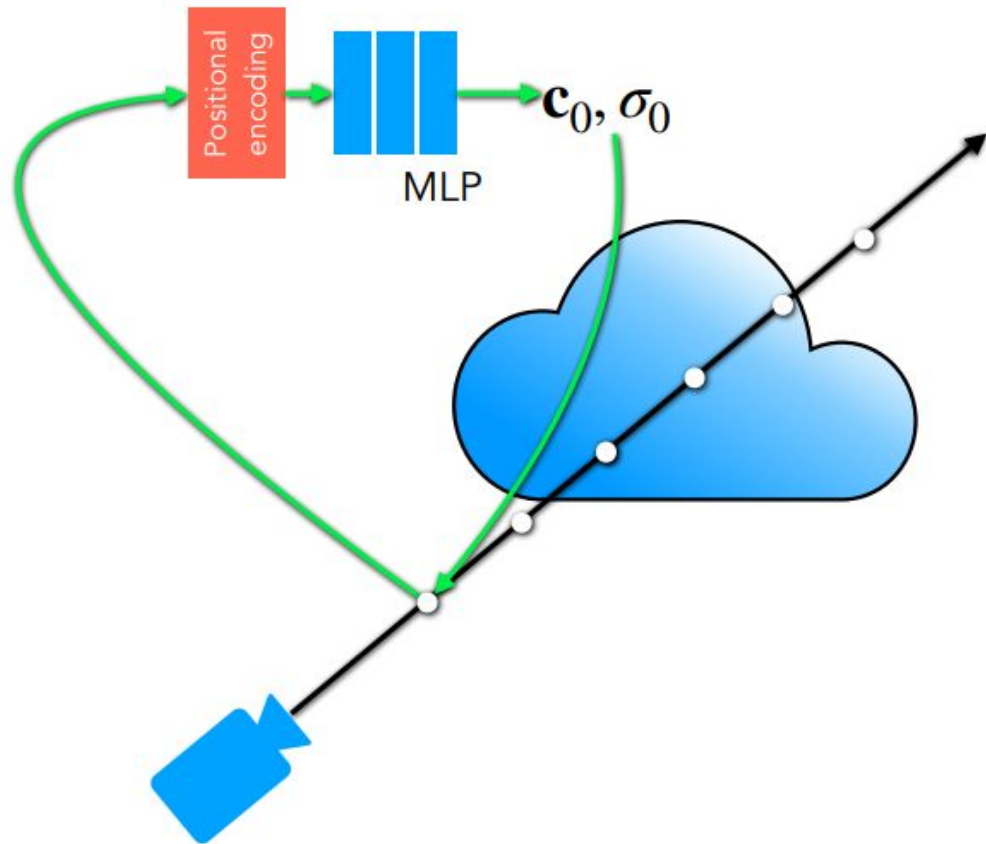
**Phew!!! Let's put everything  
together!**

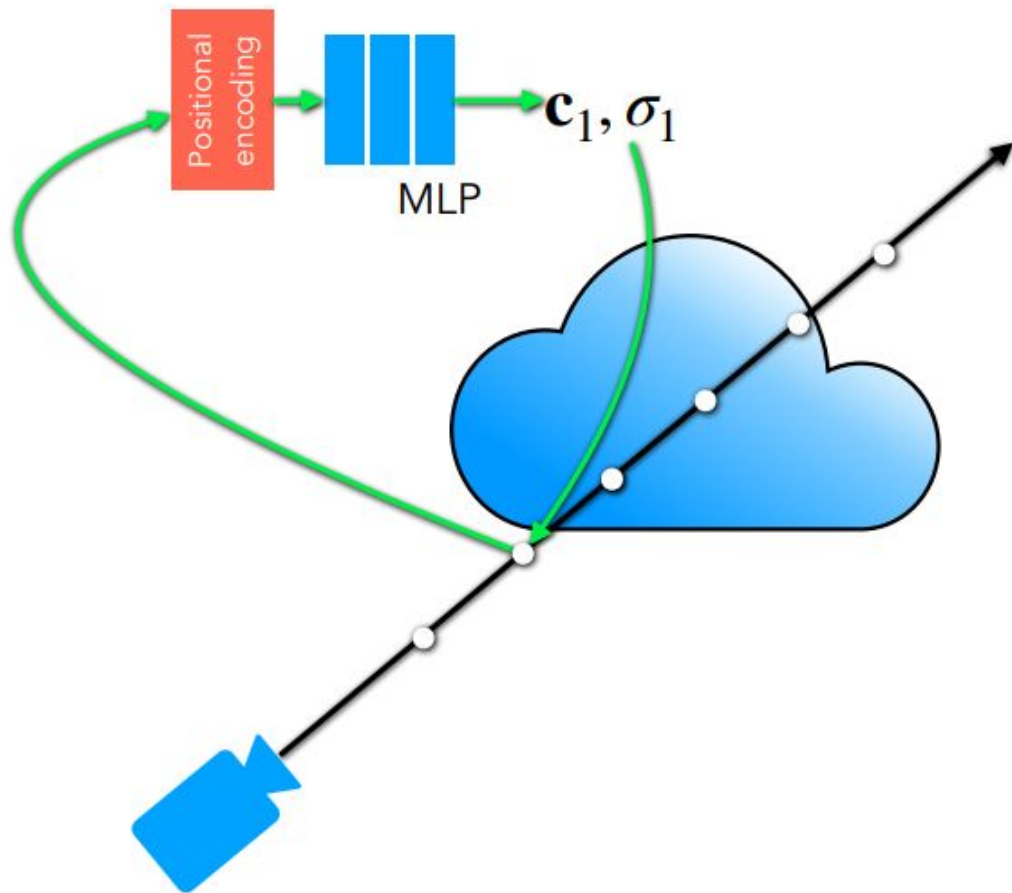
# Putting everything together!



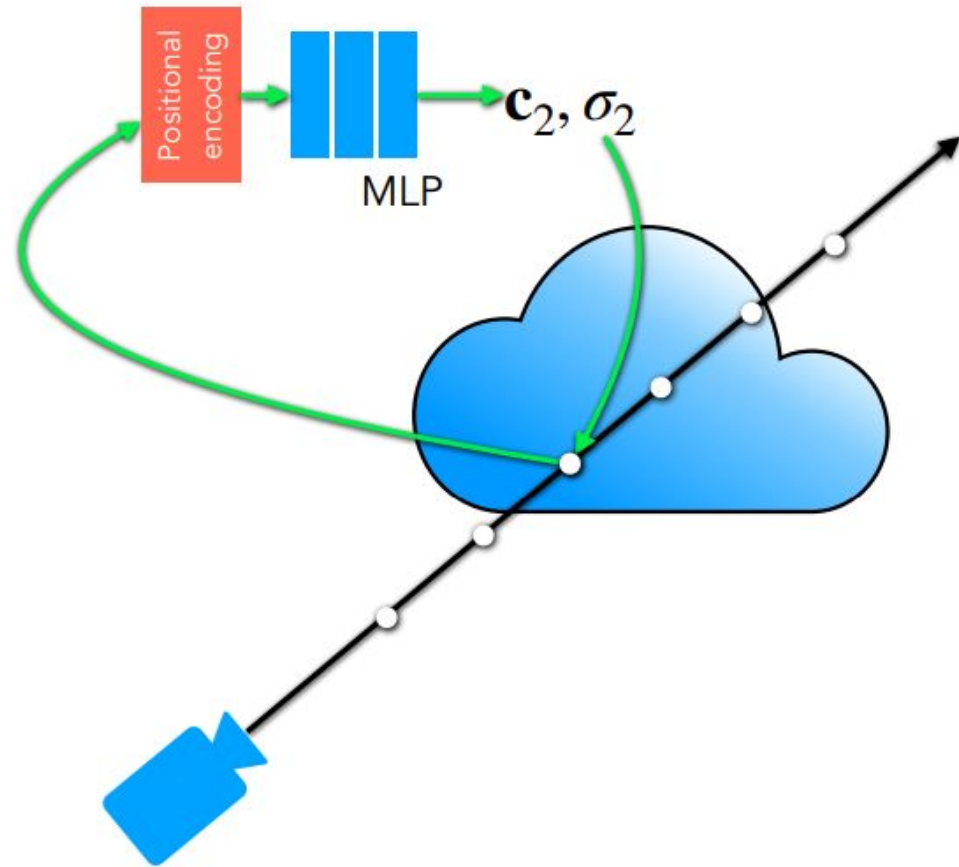


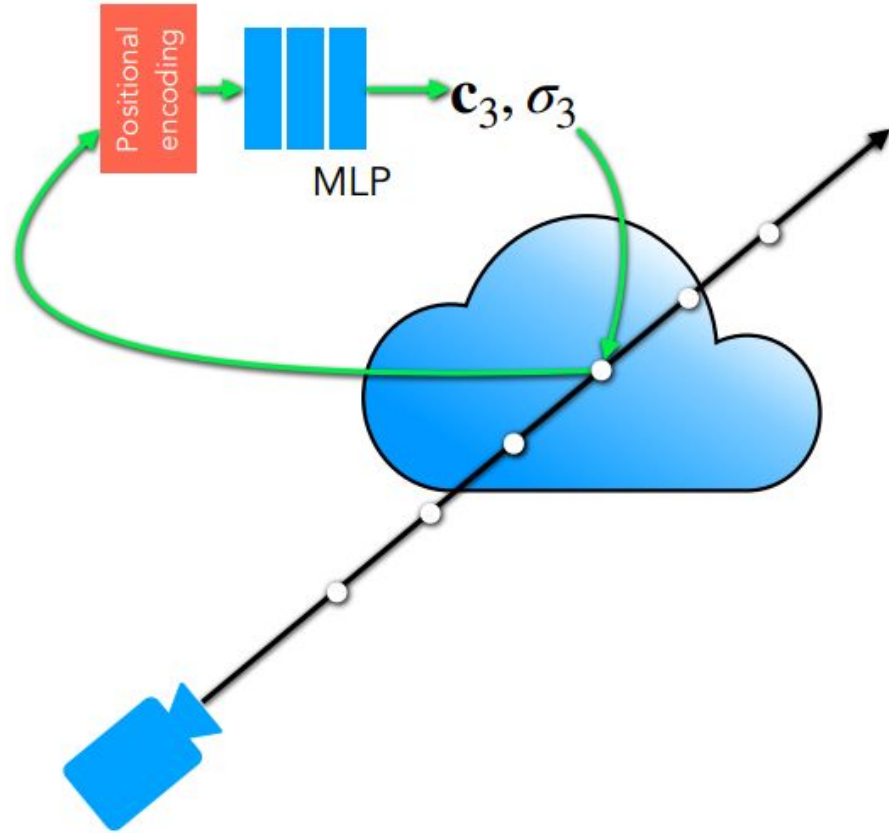


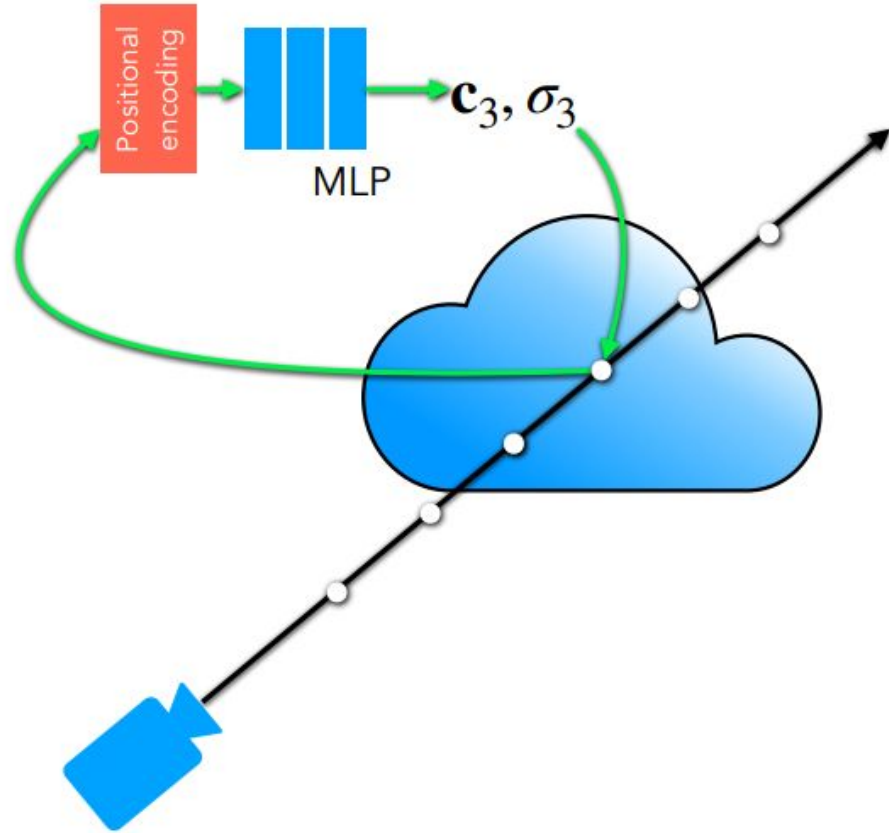


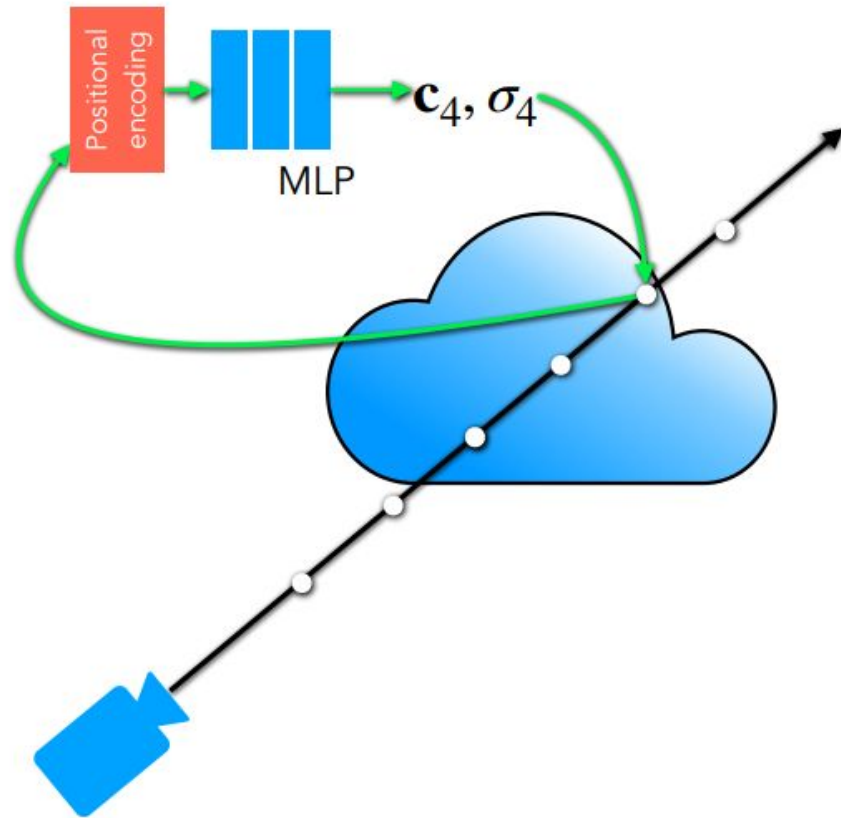


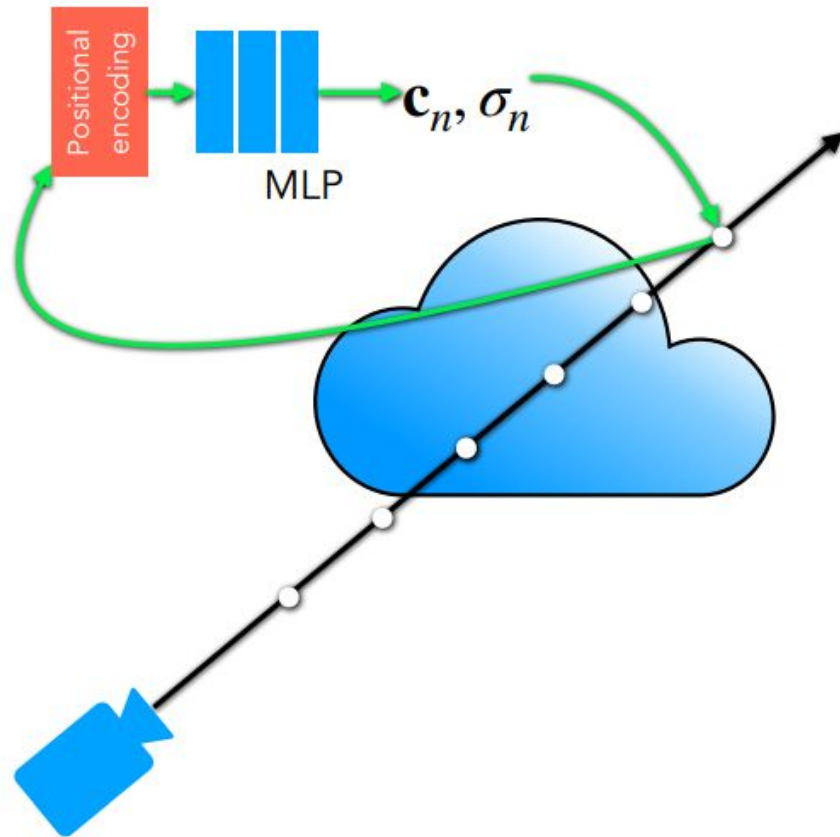














Questions?



**It's DATA TIME!!!**

# Datasets

Diffuse Synthetic 360° - 1.54 GB

- 8 Objects



Realistic Synthetic 360° - 1.56 GB

- Pinecone Images
- Flower Vase Images






# Data

Dataset	Real or Synthetic Objects	Image Size	Training Images	Testing Images
LLFF	Real	1008 × 756	20 - 62	1/8 of Training Images
Diffuse Synthetic 360	Synthetic	512 × 512	479	1000
Realistic Synthetic 360	Synthetic	512 × 512	479	1000
DeepVoxels	Synthetic	512 × 512	479	1000



# NeRF Computation/Memory Requirements

- Memory Requirements for Network Weights: 5 MB
  - Training Computation Time: ~15 hours (after 200k iterations)
  - Testing Computation Time: between 30 sec and 1 minute
  - Training computation time can vary depending on the resolution
- 



**Let's talk training!**

# Training

Generate random seed [if not provided]

```
if args.random_seed is not None:
    print('Fixing random seed', args.random_seed)
    np.random.seed(args.random_seed)
    tf.compat.v1.set_random_seed(args.random_seed)
```

# Training

Split data into sets ("Train", "Validation", "Testing")

```
if args.dataset_type == 'llff':
    images, poses, bds, render_poses, i_test = load_llff_data(args.datadir, args.factor,
                                                            recenter=True)

elif args.dataset_type == 'blender':
    .....

else:
    .....

i_train, i_val, i_test = i_split
```

```
41 def load_blender_data(basedir, half_res=False, testskip=1):
42     splits = ['train', 'val', 'test']
43     metas = {}
44     for s in splits:
45         with open(os.path.join(basedir, 'transforms_{}.json'.format(s)), 'r') as fp:
46             metas[s] = json.load(fp)
47
48     all_imgs = []
49     all_poses = []
50     counts = [0]
51     for s in splits:
52         meta = metas[s]
53         imgs = []
54         poses = []
55         if s=='train' or testskip==0:
56             skip = 1
57         else:
58             skip = testskip
59
60         for frame in meta['frames'][::skip]:
61             fname = os.path.join(basedir, frame['file_path'] + '.png')
62             imgs.append(imageio.imread(fname))
63             poses.append(np.array(frame['transform_matrix']))
64         imgs = (np.array(imgs) / 255.).astype(np.float32) # keep all 4 channels (RGBA)
65         poses = np.array(poses).astype(np.float32)
66         counts.append(counts[-1] + imgs.shape[0])
67         all_imgs.append(imgs)
68         all_poses.append(poses)
69
70     i_split = [np.arange(counts[i], counts[i+1]) for i in range(3)]
71
72     imgs = np.concatenate(all_imgs, 0)
73     poses = np.concatenate(all_poses, 0)
74
75     H, W = imgs[0].shape[:2]
76     camera_angle_x = float(meta['camera_angle_x'])
77     focal = .5 * W / np.tan(.5 * camera_angle_x)
78
79     render_poses = tf.stack([pose_spherical(angle, -30.0, 4.0) for angle in np.linspace(-180,180,40+1)[:40]], 0)
80
81     if half_res:
82         imgs = tf.image.resize_area(imgs, [400, 400]).numpy()
83         H = H//2
84         W = W//2
85         focal = focal/2.
86
87     return imgs, poses, render_poses, [H, W, focal], i_split
88
```

# Training

## Create Adam Optimizer Object

Initial Learning rate:  $5 \times 10^{-4}$   
{exponentially decays to  $5 \times 10^{-5}$ }  
Other hyperparameters set to  
Adam default ( $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  
and  $\epsilon = 10^{-7}$ )

```
# Create optimizer
lrate = args.lrate
if args.lrate_decay > 0:
    lrate =
tf.keras.optimizers.schedules.ExponentialDecay(lrate,
decay_steps=args.lrate_decay * 1000, decay_rate=0.1)
optimizer = tf.keras.optimizers.Adam(lrate)
models['optimizer'] = optimizer
```

# Training

Conduct training for N iterations. At each step:

- a. Sample random data
- b. Make predictions for parameters
- c. Computer loss (MSE)
- d. Add in the loss for the coarse grained model
- e. Apply the gradients

<https://github.com/bmild/nerf>

```
N_iters = 1000000
for i in range(start, N_iters):
    # Random from one image
    img_i = np.random.choice(i_train)
    target = images[img_i]

    # Make predictions for color, disparity, accumulated opacity.
    rgb, disp, acc, extras = render(
        H, W, focal, chunk=args.chunk, rays=batch_rays)
    # Compute MSE loss between predicted and true RGB
    img_loss = img2mse(rgb, target_s)
    loss = img_loss
    psnr = mse2psnr(img_loss)

    # Add MSE loss for coarse-grained model
    if 'rgb0' in extras:
        img_loss0 = img2mse(extras['rgb0'], target_s)
        loss += img_loss0
        psnr0 = mse2psnr(img_loss0)

    optimizer.apply_gradients(zip(gradients, grad_vars))
```

# Loss Function

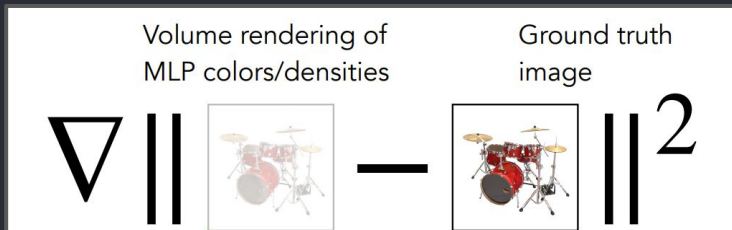
The total squared error between the rendered and true pixel colors for both the coarse and fine renderings:

$$\mathcal{L} = \sum_{\mathbf{r} \in \mathcal{R}} \left[ \left\| \hat{C}_c(\mathbf{r}) - C(\mathbf{r}) \right\|_2^2 + \left\| \hat{C}_f(\mathbf{r}) - C(\mathbf{r}) \right\|_2^2 \right]$$

```
def img2mse(x, y): return tf.reduce_mean(tf.square(x - y))
```

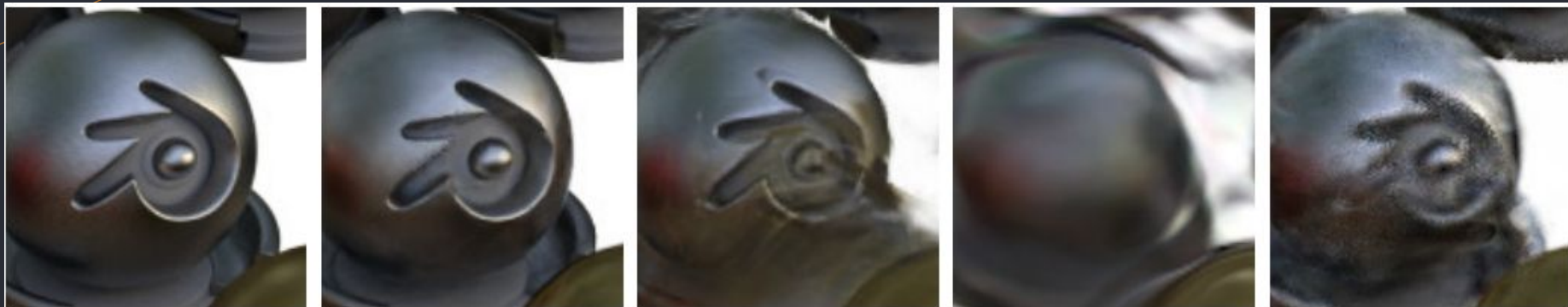
<https://arxiv.org/pdf/2003.08934.pdf>

<http://graphics.stanford.edu/courses/cs348n-22-winter/>





# Visual Results



Ground Truth

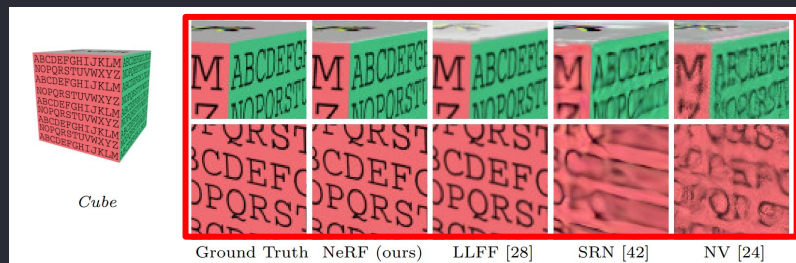
NeRF (ours)

LLFF [28]

SRN [42]

NV [24]

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



Cube

Ground Truth

NeRF (ours)

LLFF [28]

SRN [42]

NV [24]

# Comparative Results

Method	Diffuse Synthetic 360° [41]			Realistic Synthetic 360°			Real Forward-Facing [28]		
	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓
SRN [42]	33.20	0.963	0.073	22.26	0.846	0.170	22.84	0.668	0.378
NV [24]	29.62	0.929	0.099	26.05	0.893	0.160	-	-	-
LLFF [28]	34.38	0.985	0.048	24.88	0.911	0.114	24.13	0.798	<b>0.212</b>
Ours	<b>40.15</b>	<b>0.991</b>	<b>0.023</b>	<b>31.01</b>	<b>0.947</b>	<b>0.081</b>	<b>26.50</b>	<b>0.811</b>	0.250

$$\text{SSIM}(x, y) = \frac{(2\mu_x\mu_y + c_1)(2\sigma_{xy} + c_2)}{(\mu_x^2 + \mu_y^2 + c_1)(\sigma_x^2 + \sigma_y^2 + c_2)}$$

with:

- $\mu_x$  the pixel sample mean of  $x$ ;
- $\mu_y$  the pixel sample mean of  $y$ ;
- $\sigma_x^2$  the variance of  $x$ ;
- $\sigma_y^2$  the variance of  $y$ ;
- $\sigma_{xy}$  the covariance of  $x$  and  $y$ ;
- $c_1 = (k_1L)^2$ ,  $c_2 = (k_2L)^2$  two variables to stabilize the division with weak denominator;
- $L$  the dynamic range of the pixel-values (typically this is  $2^{\text{\#bits per pixel}} - 1$ );
- $k_1 = 0.01$  and  $k_2 = 0.03$  by default.

$$\text{MSE} = \frac{1}{c * i * j} \sum (I_1 - I_2)^2$$

$$\text{PSNR} = 10 \cdot \log_{10} \left( \frac{\text{MAX}_I^2}{\text{MSE}} \right)$$

LPIPS:  
<https://arxiv.org/pdf/1801.03924.pdf>

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

# Ablation Study Results

	Input	#Im.	$L$	$(N_c, N_f)$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
1) No PE, VD, H	$xyz$	100	-	(256, -)	26.67	0.906	0.136
2) No Pos. Encoding	$xyz\theta\phi$	100	-	(64, 128)	28.77	0.924	0.108
3) No View Dependence	$xyz$	100	10	(64, 128)	27.66	0.925	0.117
4) No Hierarchical	$xyz\theta\phi$	100	10	(256, -)	30.06	0.938	0.109
5) Far Fewer Images	$xyz\theta\phi$	25	10	(64, 128)	27.78	0.925	0.107
6) Fewer Images	$xyz\theta\phi$	50	10	(64, 128)	29.79	0.940	0.096
7) Fewer Frequencies	$xyz\theta\phi$	100	5	(64, 128)	30.59	0.944	0.088
8) More Frequencies	$xyz\theta\phi$	100	15	(64, 128)	30.81	0.946	0.096
9) Complete Model	$xyz\theta\phi$	100	10	(64, 128)	<b>31.01</b>	<b>0.947</b>	<b>0.081</b>

Table 2: An ablation study of our model. Metrics are averaged over the 8 scenes from our realistic synthetic dataset. See Sec. 6.4 for detailed descriptions.

# Summary

- ❖ NeRF = Coordinate Based MLP Network + Volume Rendering
- ❖ NeRF Pros:
  - Simple representation
  - Differentiable Rendering Model
- ❖ NeRF Cons:
  - Dumb Brute force approach
  - INSANELY SLOW!!!



Original NeRF

Training speed

1-2 days

Rendering speed

30 sec



KiloNeRF,  
cached voxels

1-2 days

1/60 sec



Learned voxels

10-15 mins

1/15-1/2 sec



Learned hash maps  
(Instant NGP)

5 sec - 5 mins

1/60 sec



Questions?

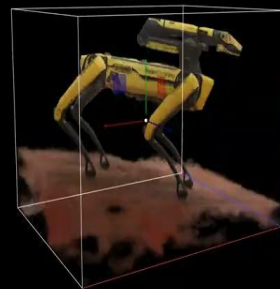
# DeepRob

[Student] Lecture 17

by *Michael Andrev, Chen Hu, Sean Coffey*

Implicit Surfaces, Geometry, NeRF

University of Michigan and University of Minnesota



# Sean References

Moons, Theo, Luc Van Gool, and Maarten Vergauwen. "3D reconstruction from multiple images part 1: Principles." Foundations and Trends in Computer Graphics and Vision 4.4 (2010): 287-404.

Mahmoudzadeh, Ahmadreza; Golroo, Amir; Jahanshahi, Mohammad R.; Firoozi Yeganeh, Sayna (January 2019). "Estimating Pavement Roughness by Fusing Color and Depth Data Obtained from an Inexpensive RGB-D Sensor". Sensors. 19 (7): 1655. Bibcode:2019Senso..19.1655M. doi:10.3390/s19071655. PMC 6479490. PMID 30959936

Buelthoff, Heinrich H., and Alan L. Yuille. "Shape-from-X: Psychophysics and computation Archived 2011-01-07 at the Wayback Machine." Fibers '91, Boston, MA. International Society for Optics and Photonics, 1991

Richard Hartley and Andrew Zisserman (2003). "Multiple View Geometry in computer vision". Cambridge University Press. ISBN 978-0-521-54051-3

Yuan Yao, Yasamin Jafarian, & Hyun Soo Park. (2019). "MONET: Multiview Semi-supervised Keypoint Detection via Epipolar Divergence." arXiv:1806.00104. <https://arxiv.org/pdf/1806.00104.pdf>