

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

조민지

Vision & Display Systems Lab.

Dept. of Electronic Engineering, Sogang University

Outline

- Intro
- Previous works
- NeRF : Representing Scenes as Neural Radiance Fields for View Synthesis
- Method
- Results
- Going forward

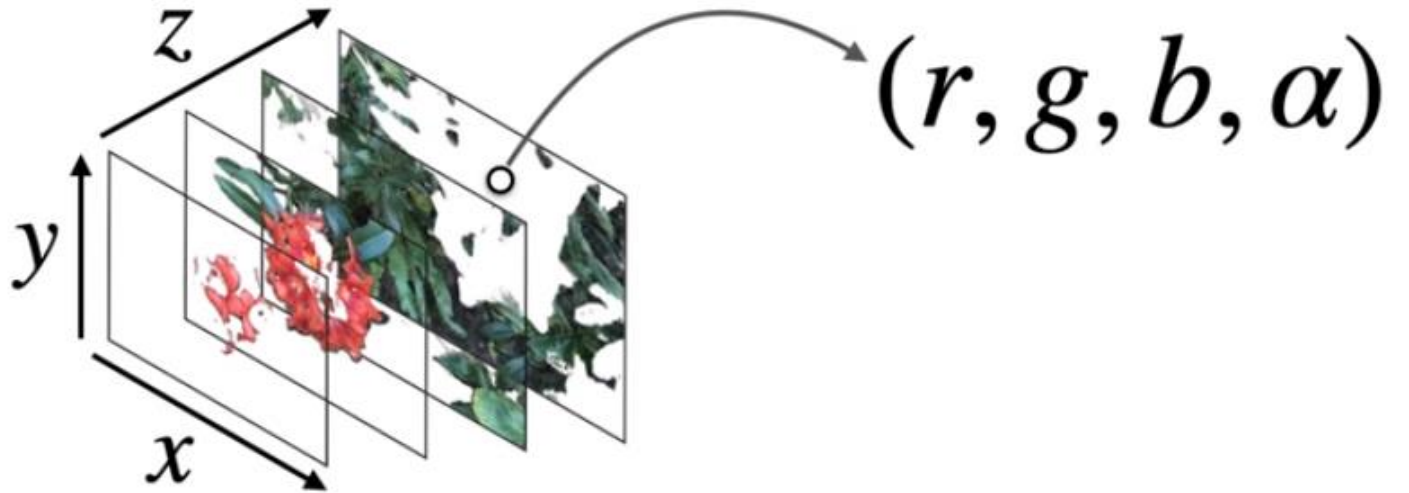
Intro

- 3D Rendering
and View synthesis



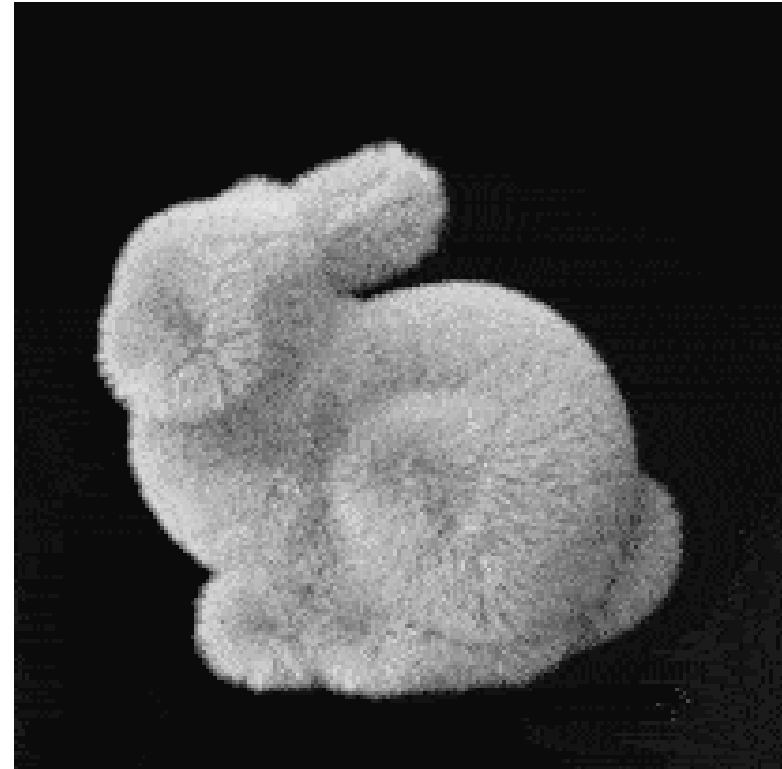
Previous works

- RGB- α volume rendering for synthesis
 - Large N-d array contains RGB- α information.



Previous works

- Neural Network as shape representation
 - Deep SDF
 - Occupancy network



NeRF : Neural Radiance Field Scene Representation

- Neural Network as scene representation
 - Volume rendering with neural radiance fields
 - View synthesis and image-based rendering



Ben Mildenhall

Method

- Neural Network as a scene representation
 - Generalization? $\rightarrow X$
 - Overfitted by 1 specific scene.
 - Weights of NN represents the scene!



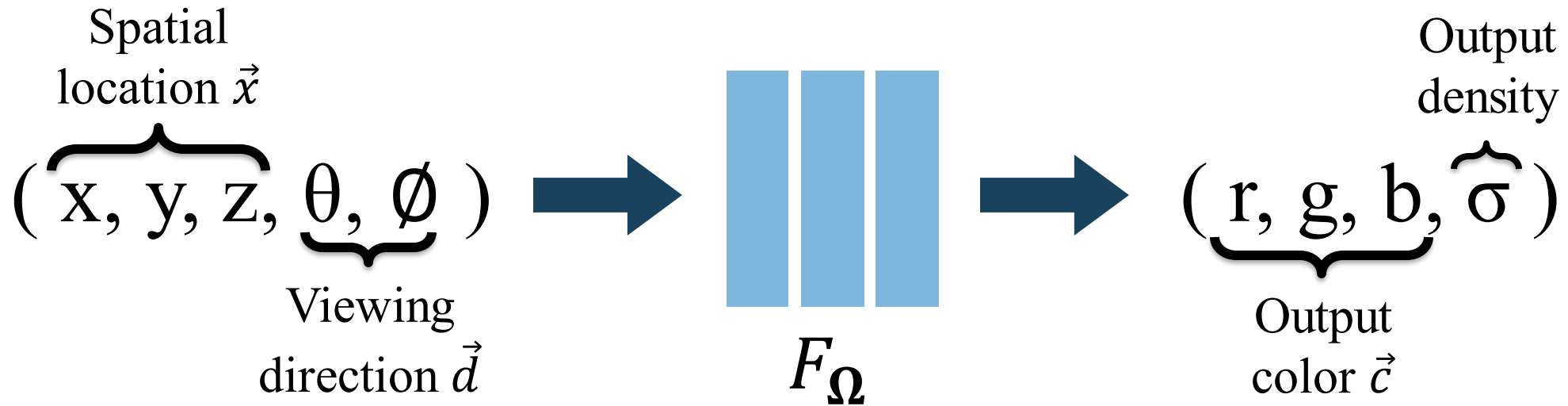
F_{Ω}

Fully-connected
neural network

(256 channels,
9 layers)

Method

- Representing a scene as a continuous 5D function



Method

- Getting pose & bound from images
 - Capture images of ¹ forward facing view or ² 360 ° inward facing view.
 - Use COLMAP to get camera-to-world (c2w) transformation matrix.



Method

- Getting pose & bound from images
 - Capture images of ¹ forward facing view or ² 360 ° inward facing view.
 - Use COLMAP to get camera-to-world (c2w) transformation matrix.



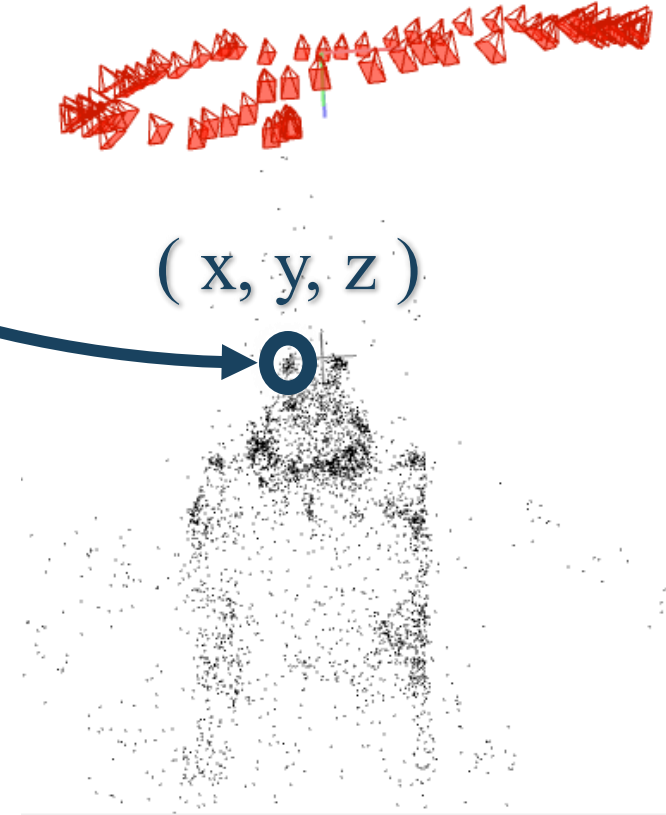
Method

- Getting pose & bound from images
 - Capture images of ¹ forward facing view or ² 360 ° inward facing view.
 - Use COLMAP to get camera-to-world (c2w) transformation matrix.



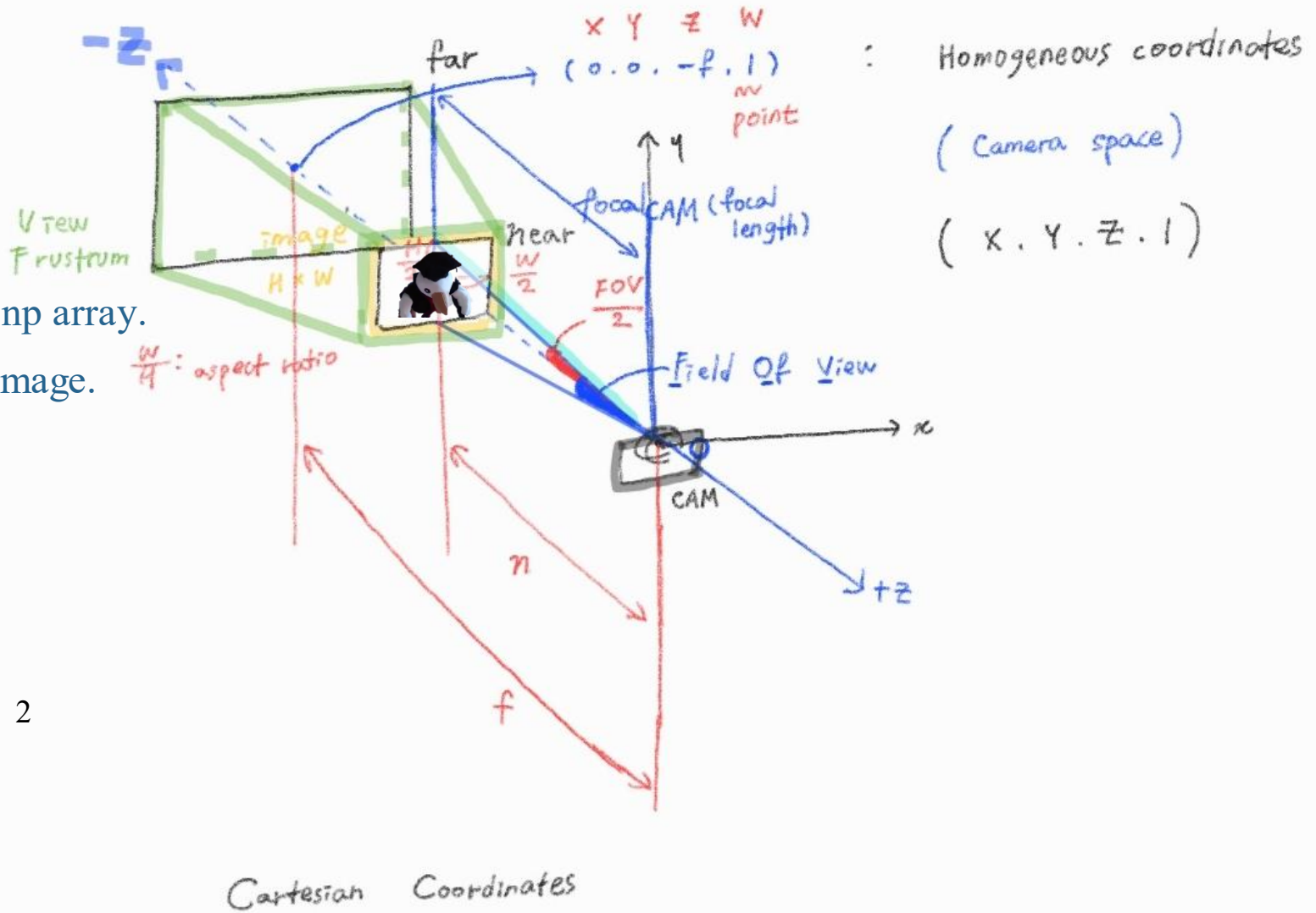
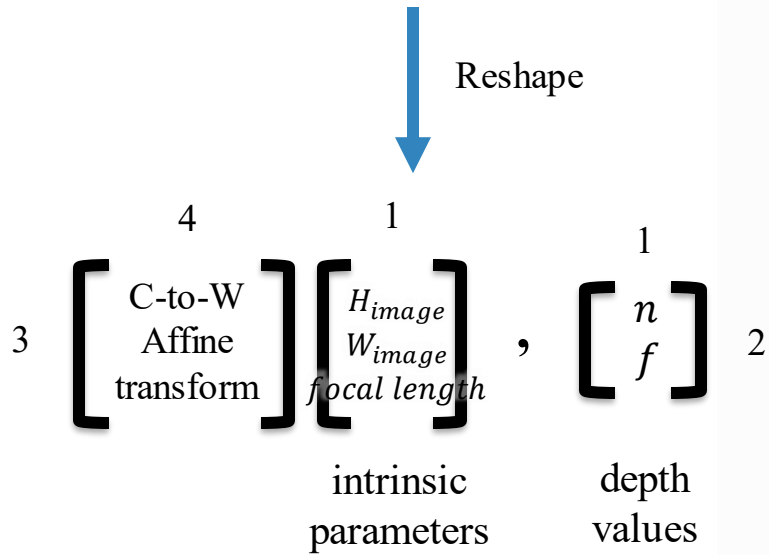
Method

- Getting pose & bound from
 - Capture images of ¹ forward or ² 360° inward facing view
 - Use COLMAP to get camera-to-world (c2w) transformation matrix.



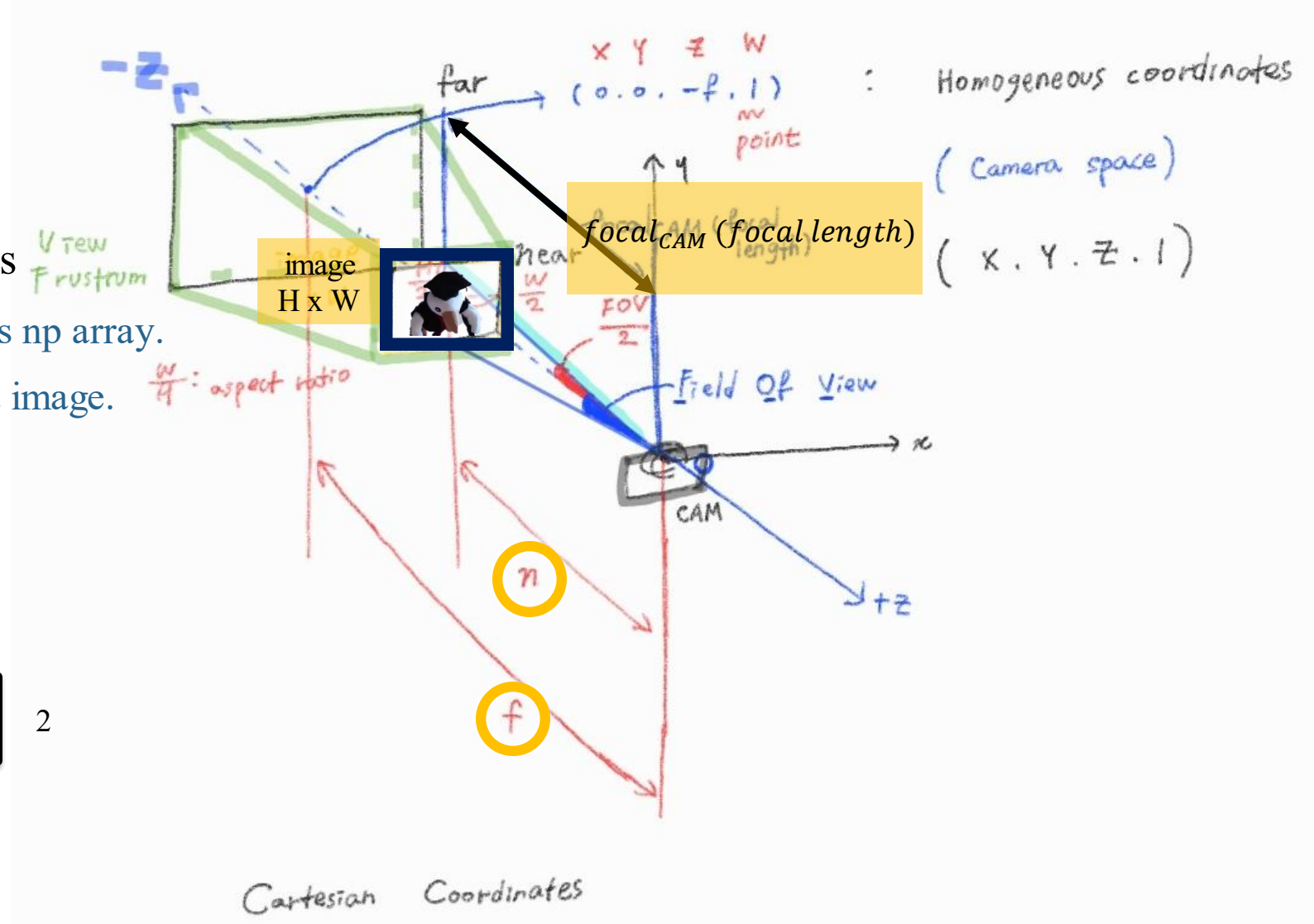
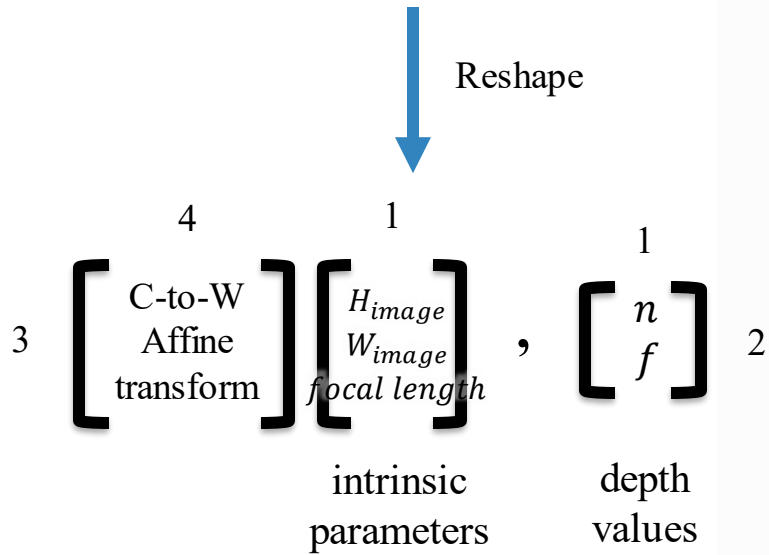
Method

- Getting pose & bound from images
 - COLMAP generates poses_bounds np array.
 - Contains 1 x 17 matrix per 1 input image.



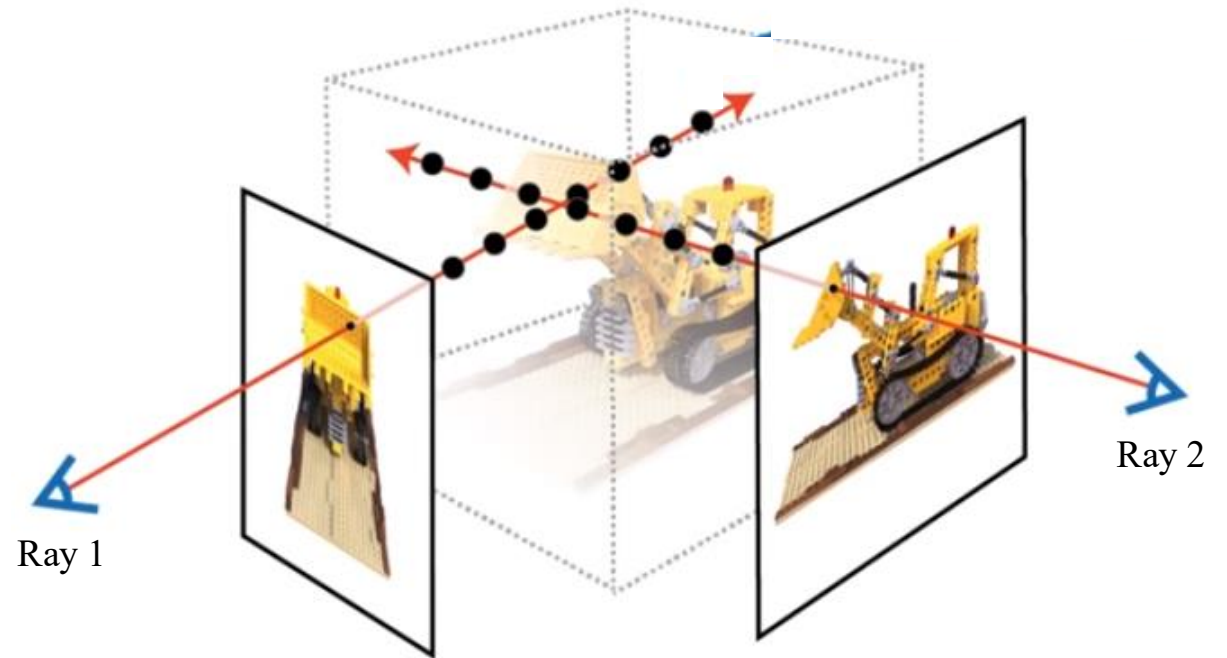
Method

- Getting pose & bound from images
 - COLMAP generates poses_bounds np array.
 - Contains 1 x 17 matrix per 1 input image.



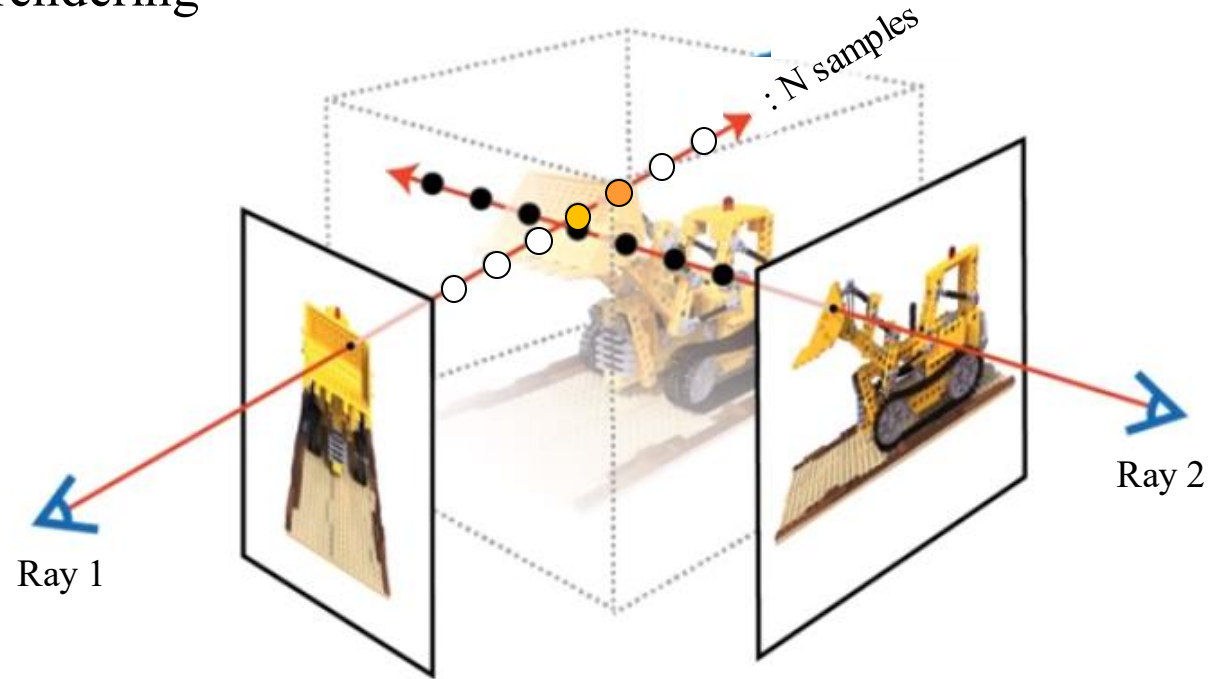
Method

- Generate views with traditional volume rendering
 - 1 pixel of each image = 1 ray



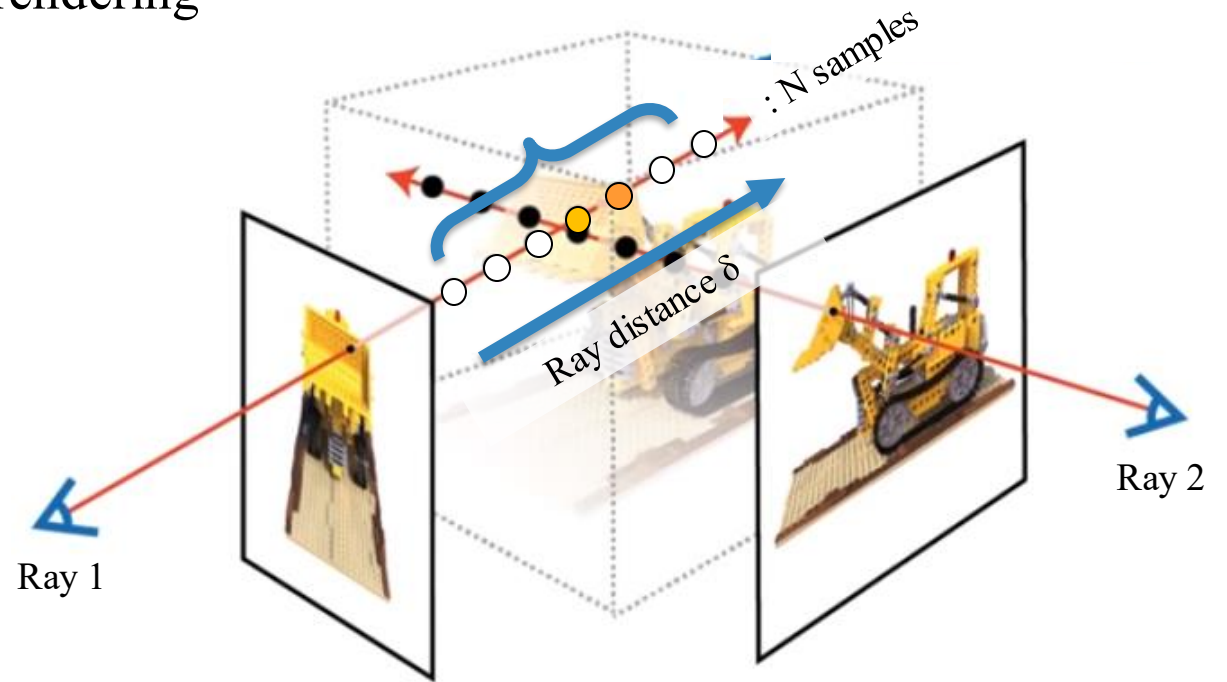
Method

- Generate views with traditional volume rendering
 - 1 pixel of each image = 1 ray



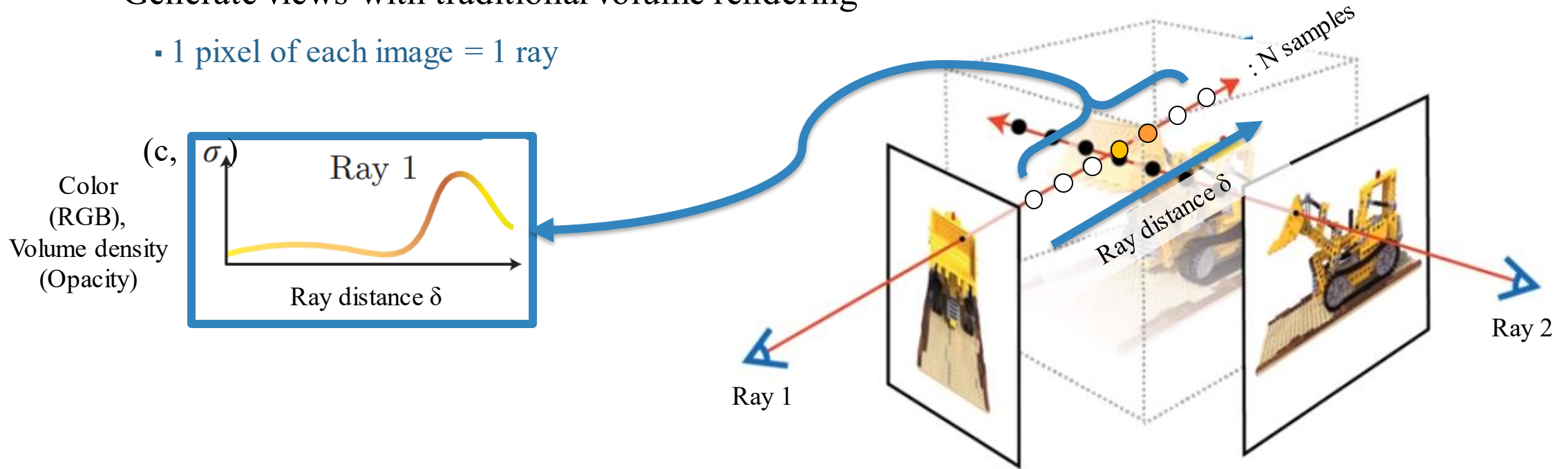
Method

- Generate views with traditional volume rendering
 - 1 pixel of each image = 1 ray



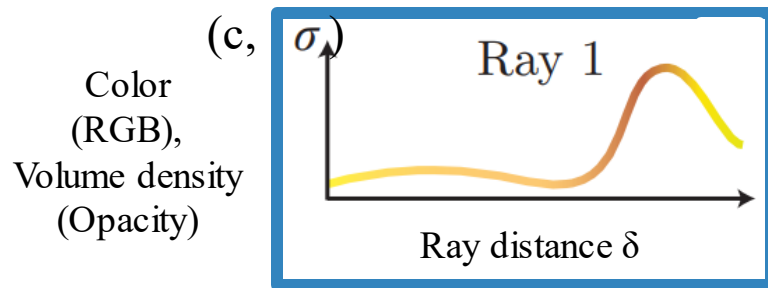
Method

- Generate views with traditional volume rendering
 - 1 pixel of each image = 1 ray



Method

- Generate views with traditional volume rendering
 - Optimize every ray with gradient descent (l2 loss)



$$\min_{\Omega} \sum_i \left\| \text{render}^{(i)}(F_{\Omega}) - I_{\text{gt}}^{(i)} \right\|^2 \quad : \mathcal{C}(i)$$

$$: \hat{\mathcal{C}}(i) = \sum_{j=1}^N T_j (1 - \exp(-\sigma_j \delta_j)) c_j$$

$$T_j = \exp(-\sum_{k=1}^{j-1} \sigma_k \delta_k)$$

Method

- Optimizing a Neural Radiance Field
 - ¹ (Sinusoidal) Positional encoding
 - For $F_{\Omega} : (x, d) \rightarrow (c, \sigma)$,
 $F_{\Omega} = F'_{\Omega} \circ \gamma$ (\circ is elementwise product)

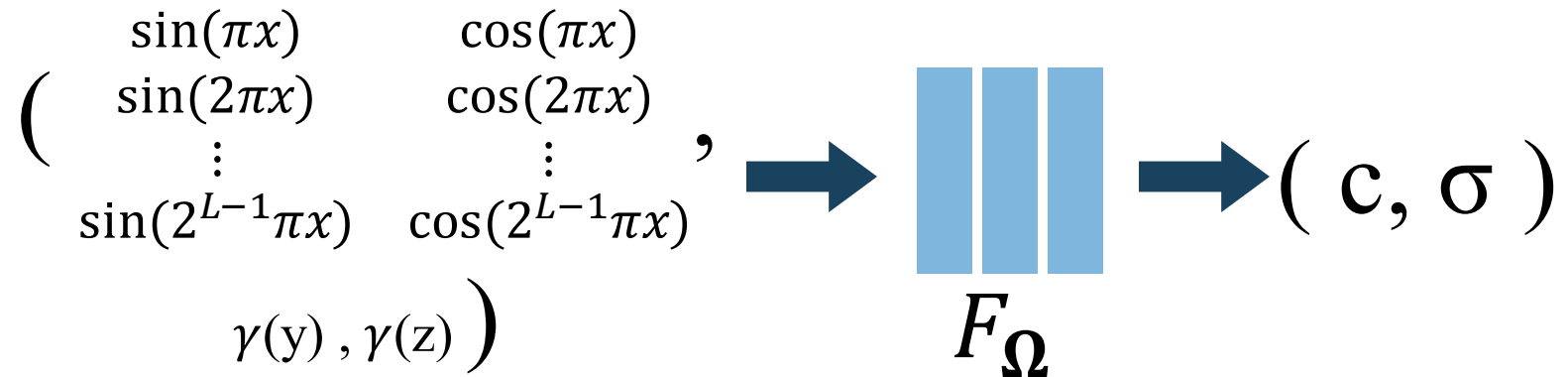
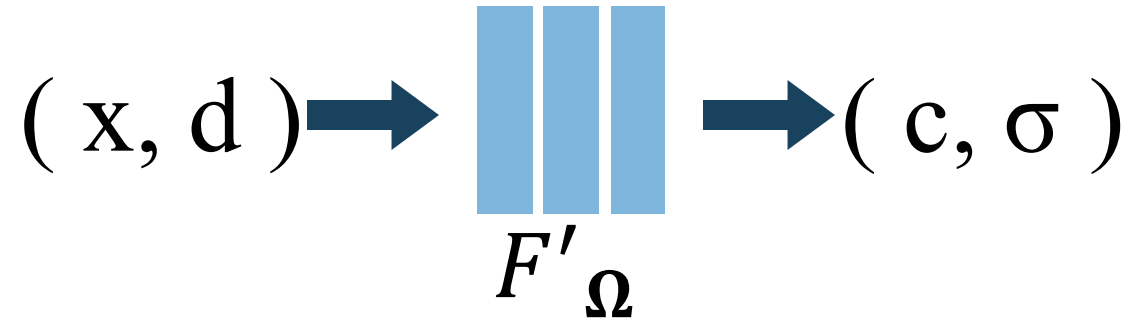
Method

- Optimizing a Neural Radiance Field

- ¹ (Sinusoidal) Positional encoding

- For $F_{\Omega} : (x, d) \rightarrow (c, \sigma)$,

- $F_{\Omega} = F'_{\Omega} \circ \gamma$ (\circ is elementwise product)



Method

- Optimizing a Neural Radiance Field

- ¹ (Sinusoidal) Positional encoding

- For $F_{\Omega} : (x, d) \rightarrow (c, \sigma)$,

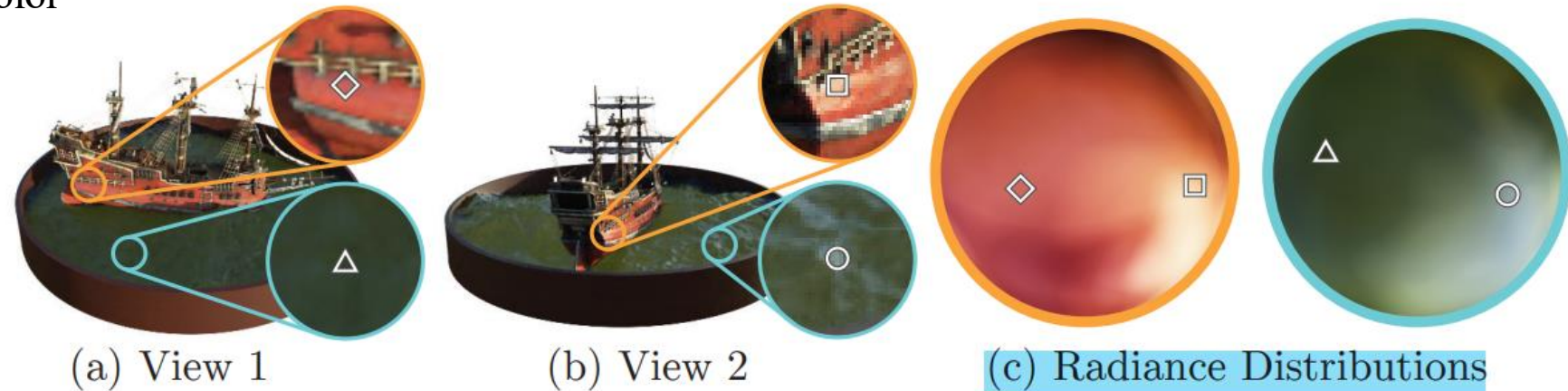
- $F_{\Omega} = F'_{\Omega} \circ \gamma$ (\circ is elementwise product)

- $\gamma(p) = (\sin(\pi p) , \cos(\pi p) , \dots , \sin(2^{L-1}\pi p) , \cos(2^{L-1}\pi p))$ ($L = 10$ for $\gamma(x)$ and 4 for $\gamma(d)$)

- $R^{2L} \leftarrow R$: map into higher dimensional space.

Method

- Optimizing a Neural Radiance Field
 - ² View-dependent RGB color



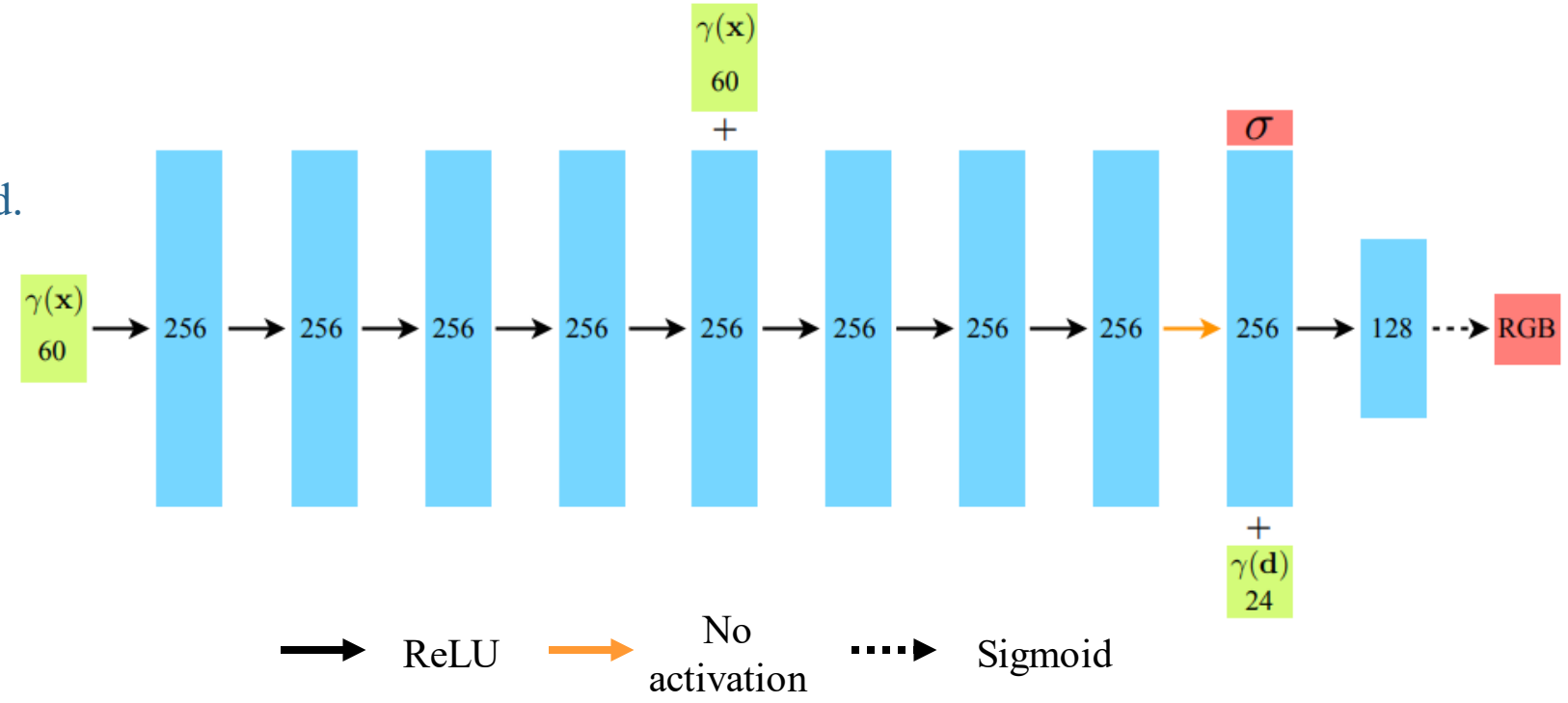
- Color (radiance) distributions of the same point on ships viewed by different angles.

Fig. 3: A visualization of view-dependent emitted radiance. Our neural radiance field representation outputs RGB color as a 5D function of both spatial position \mathbf{x} and viewing direction \mathbf{d} . Here, we visualize example **directional color distributions for two spatial locations** in our neural representation of the *Ship* scene. In (a) and (b), we show the appearance of two fixed 3D points from two different camera positions: one on the side of the ship (orange insets) and one on the surface of the water (blue insets). Our method predicts the changing specular appearance of these two 3D points, and in (c) we show how this behavior **generalizes continuously across the whole hemisphere of viewing directions**.

Method

- Optimizing a Neural Radiance Field

- ² View-dependent RGB color
- Predict σ as function of x ,
predict c as function of x and d .

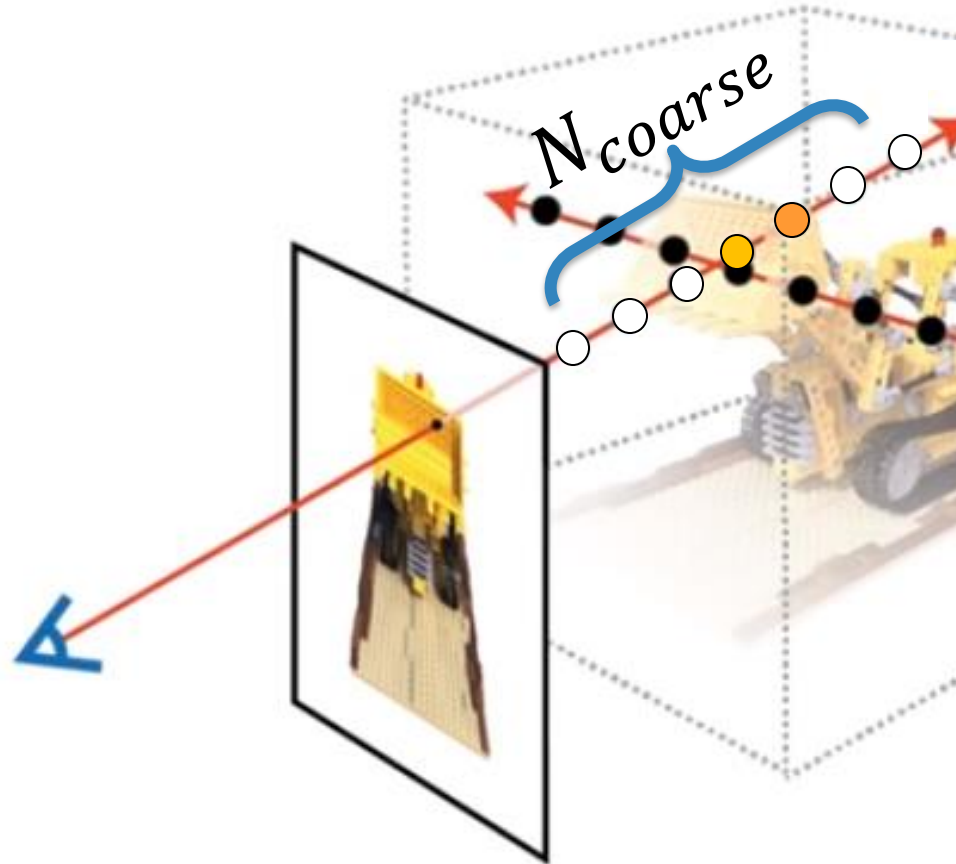


Method



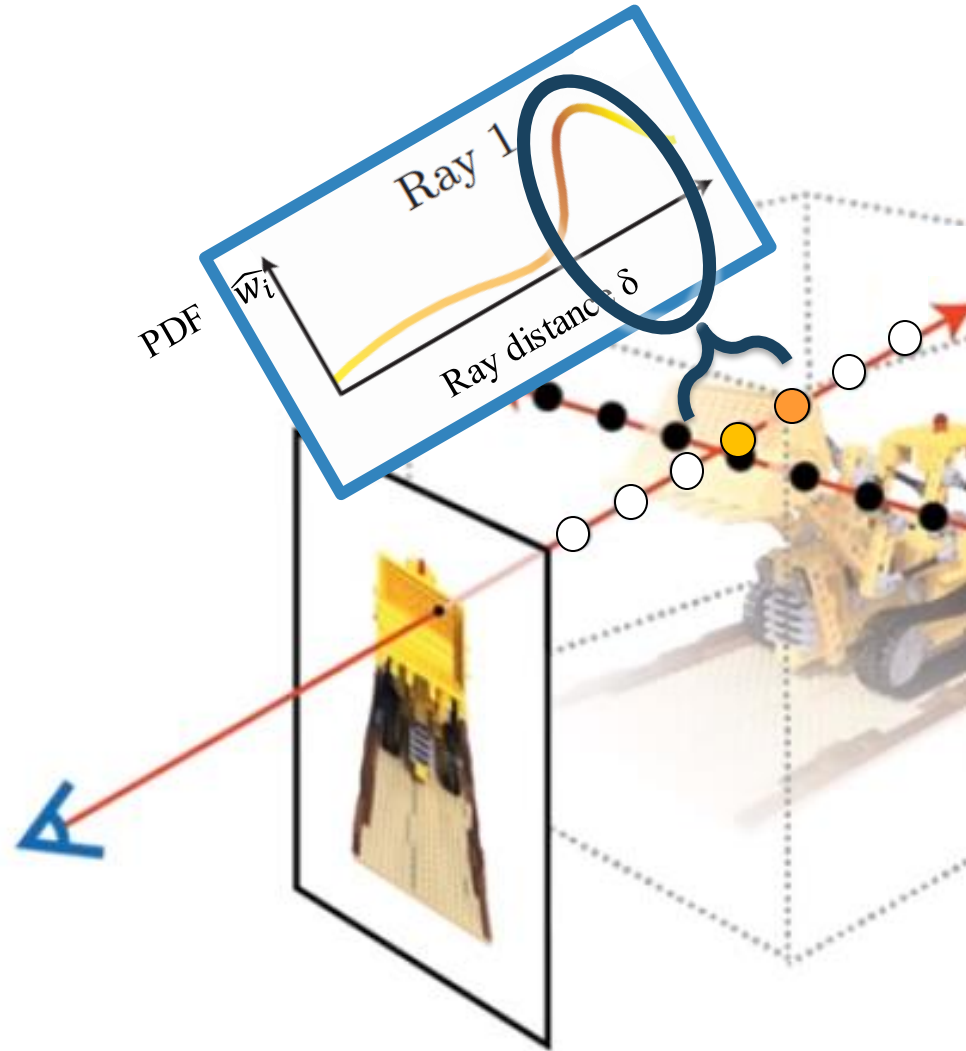
Method

- Optimizing a Neural Radiance Field
 - ³ Hierarchical volume sampling
 - Render by 2 networks : coarse & fine



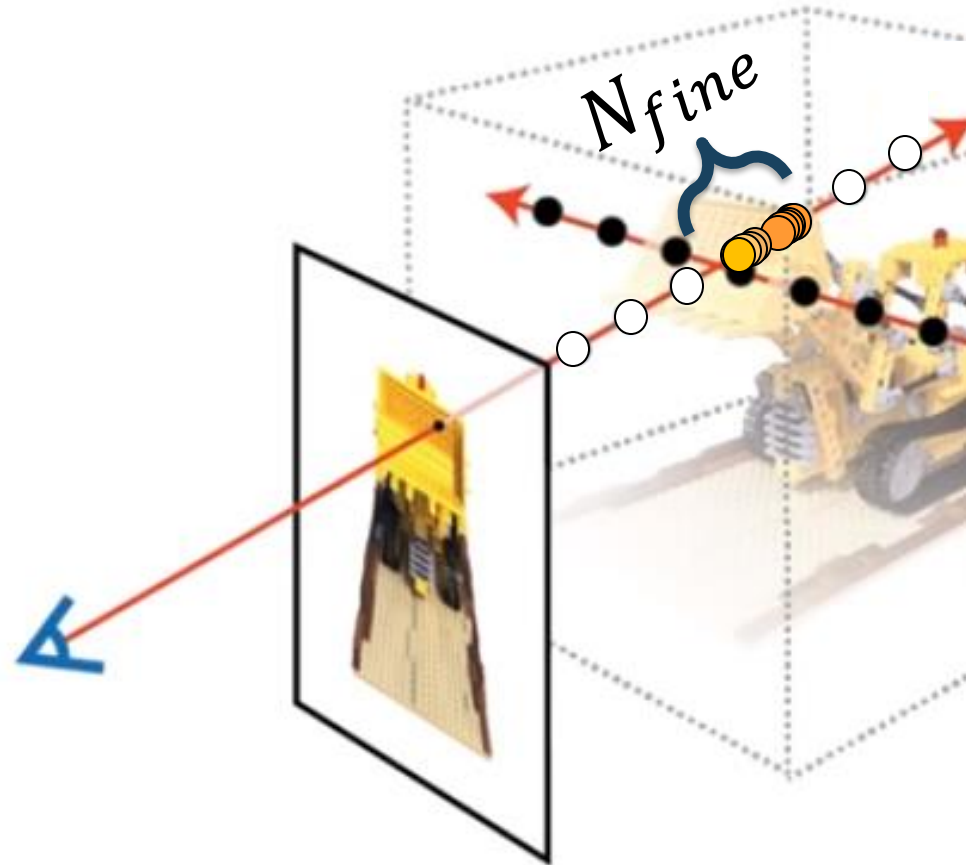
Method

- Optimizing a Neural Radiance Field
 - ³ Hierarchical volume sampling
 - Render by 2 networks : *coarse* & *fine*
 - Get Probability Distribution Function by normalizing weights from N_{coarse} .
 - Given the output of N_{coarse} , produce more informed sample points N_{fine} .
 - $N_{fine} = N_{importance}$
 - Finally, use all $N_{coarse+fine}$ samples.



Method

- Optimizing a Neural Radiance Field
 - ³ Hierarchical volume sampling
 - Render by 2 networks : *coarse* & *fine*
 - Get Probability Distribution Function by normalizing weights from N_{coarse} .
 - Given the output of N_{coarse} , produce more informed sample points N_{fine} .
 - $N_{fine} = N_{importance}$
 - Finally, use all $N_{coarse+fine}$ samples.



Results

- Ablation studies
 - Effect of 3 optimization methods

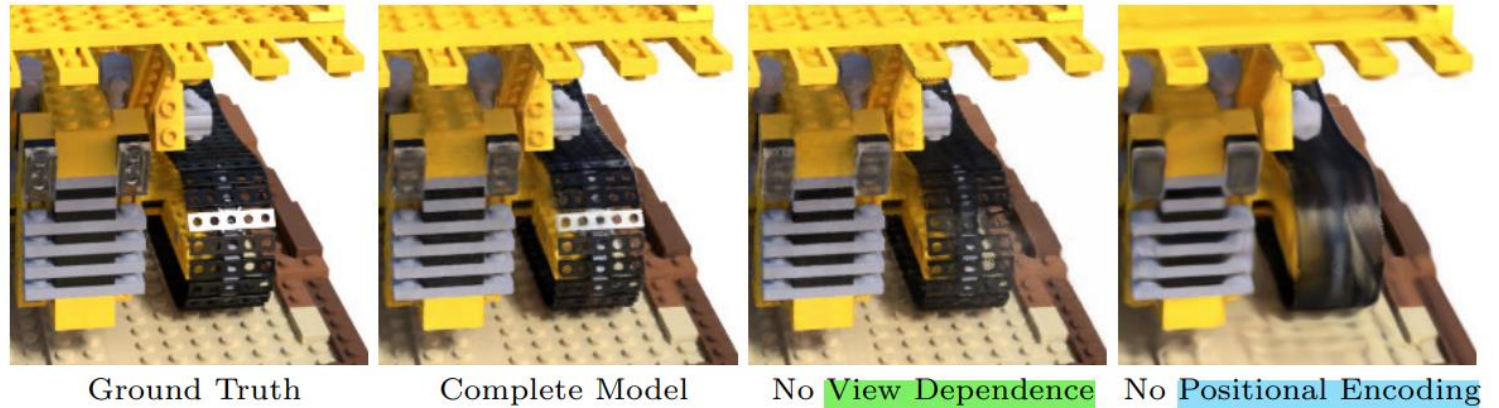
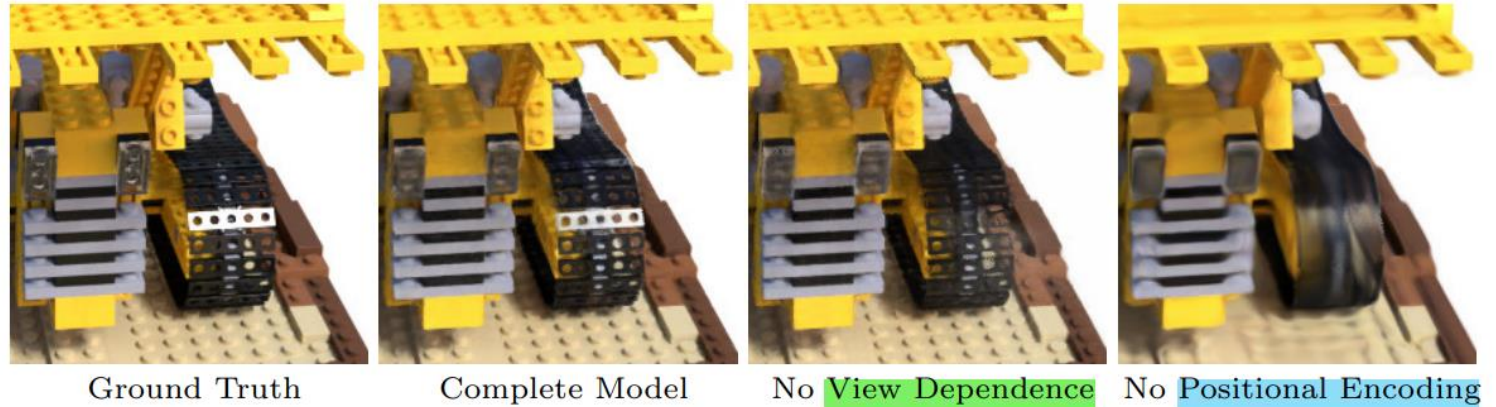


Fig. 4: Here we visualize how our full model benefits from representing view-dependent emitted radiance and from passing our input coordinates through a high-frequency positional encoding. Removing **view dependence** prevents the model from recreating the specular reflection on the bulldozer tread. Removing the **positional encoding** drastically decreases the model's ability to represent high frequency geometry and texture, resulting in an oversmoothed appearance.

Results

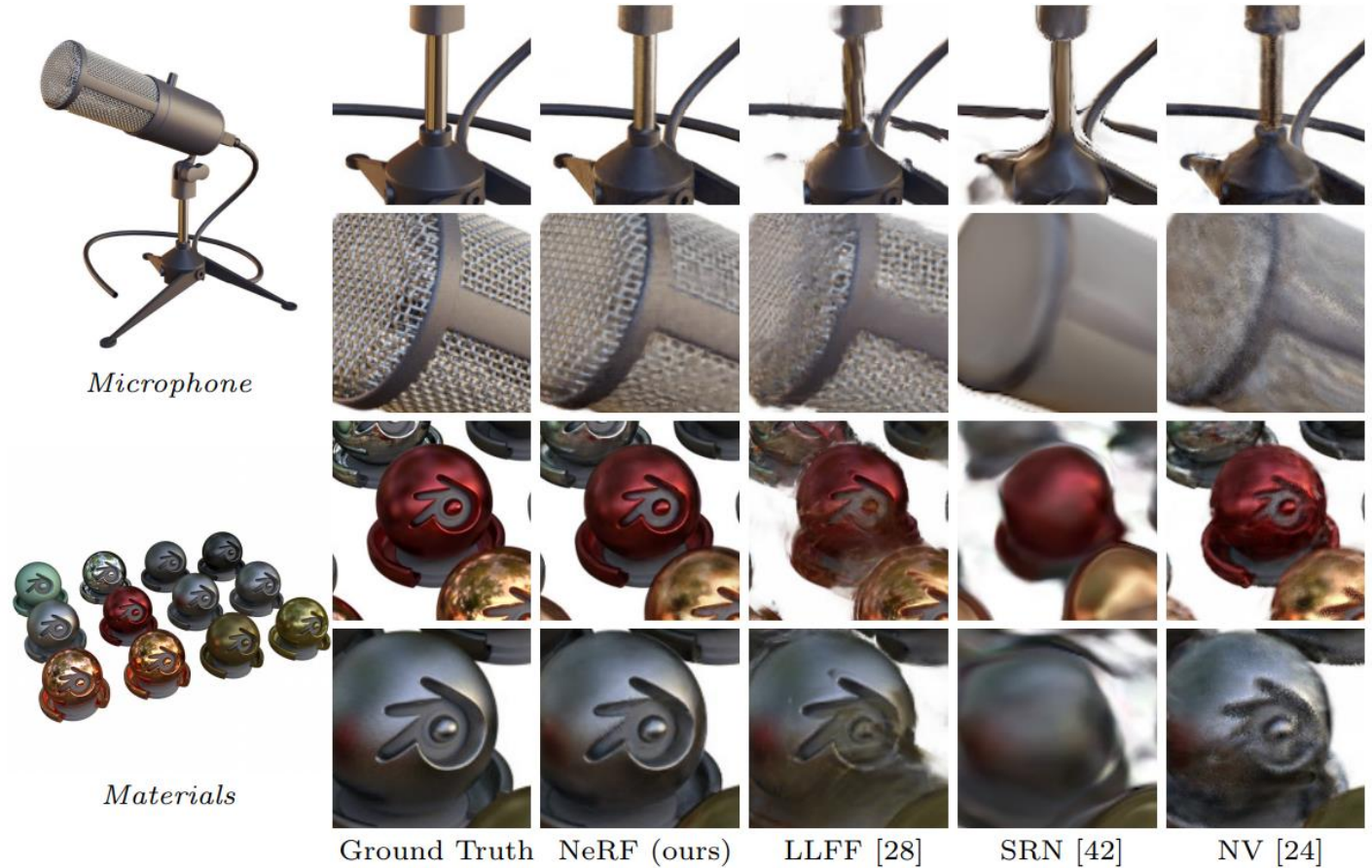
- Ablation studies
 - Effect of 3 optimization methods



	Input	#Im.	L	(N_c, N_f)	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
No PE, VD, H	xyz	100	-	(256, -)	26.67	0.906	0.136
No Pos. Encoding	$xyz\theta\phi$	100	-	(64, 128)	28.77	0.924	0.108
No View Dependence	xyz	100	10	(64, 128)	27.66	0.925	0.117
No Hierarchical	$xyz\theta\phi$	100	10	(256, -)	30.06	0.938	0.109
Complete Model	$xyz\theta\phi$	100	10	(64, 128)	31.01	0.947	0.081

Results

- Comparisons with previous methods
 - Synthetic dataset



Results

- Comparisons with previous methods
 - Real world scenes



Fern



T-Rex



Ground Truth

NeRF (ours)

LLFF [28]

SRN [42]



Results



Going forward

- Decrease training and inference time (currently 30 sec per high-resolution frame)
 - Mip-NeRF (A Multiscale Representation for Anti-Aliasing Neural Radiance Fields)
<https://jonbarron.info/mipnerf/>
- Disentangle more graphics attributes to allow additional applications, such as relighting
 - NeRV (Neural Reflectance and Visibility Fields for Relighting and View Synthesis)
<https://pratulsrinivasan.github.io/nerv/>
- Generalize to more than one scene without training from scratch and requiring fewer input images

감사합니다