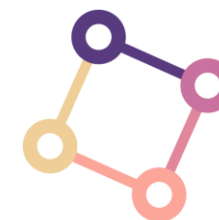


NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis

ECCV 2020 Best Paper Honorable Mention

박성현



DAVIAN

Data and Visual Analytics Lab

Novel view synthesis

Input Images



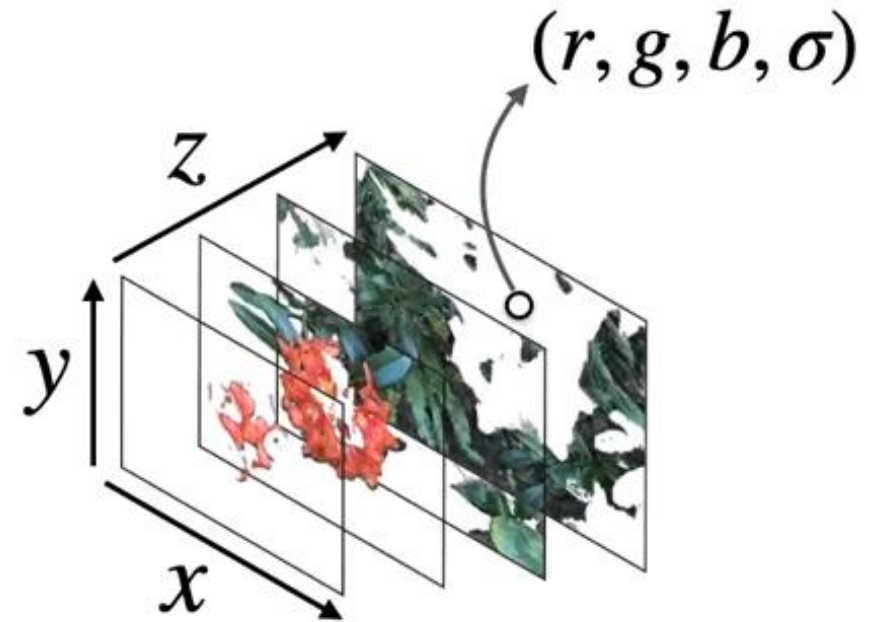
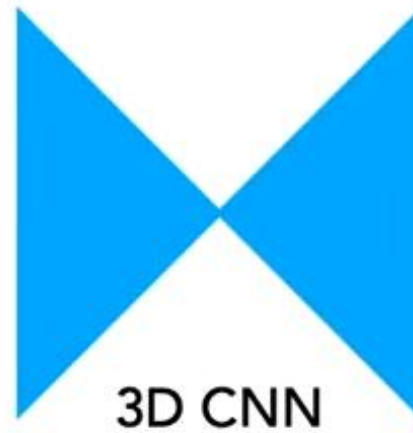
Optimize NeRF



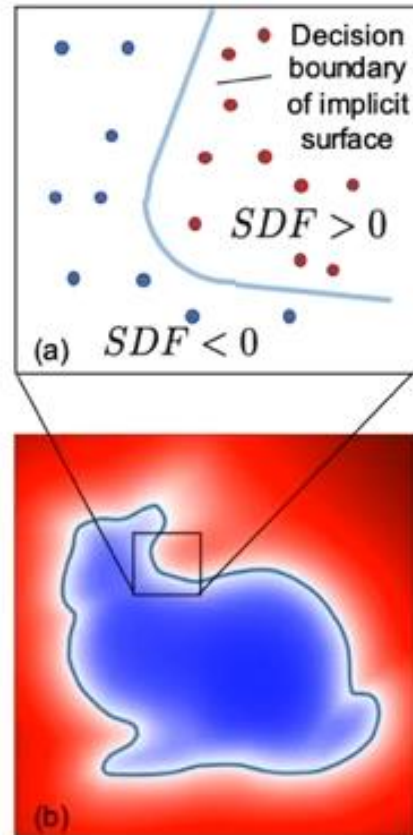
Render new views



Predict 3D voxel RGB-alpha grid



Neural networks as a shape representation



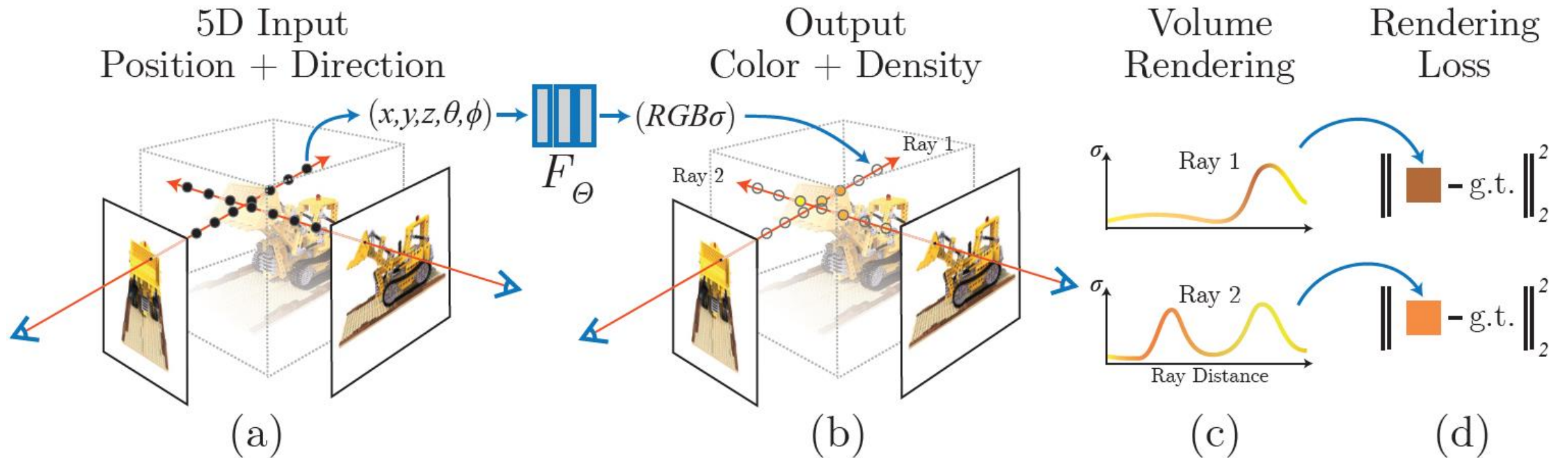
DeepSDF, Park et al. 2019



Contribution

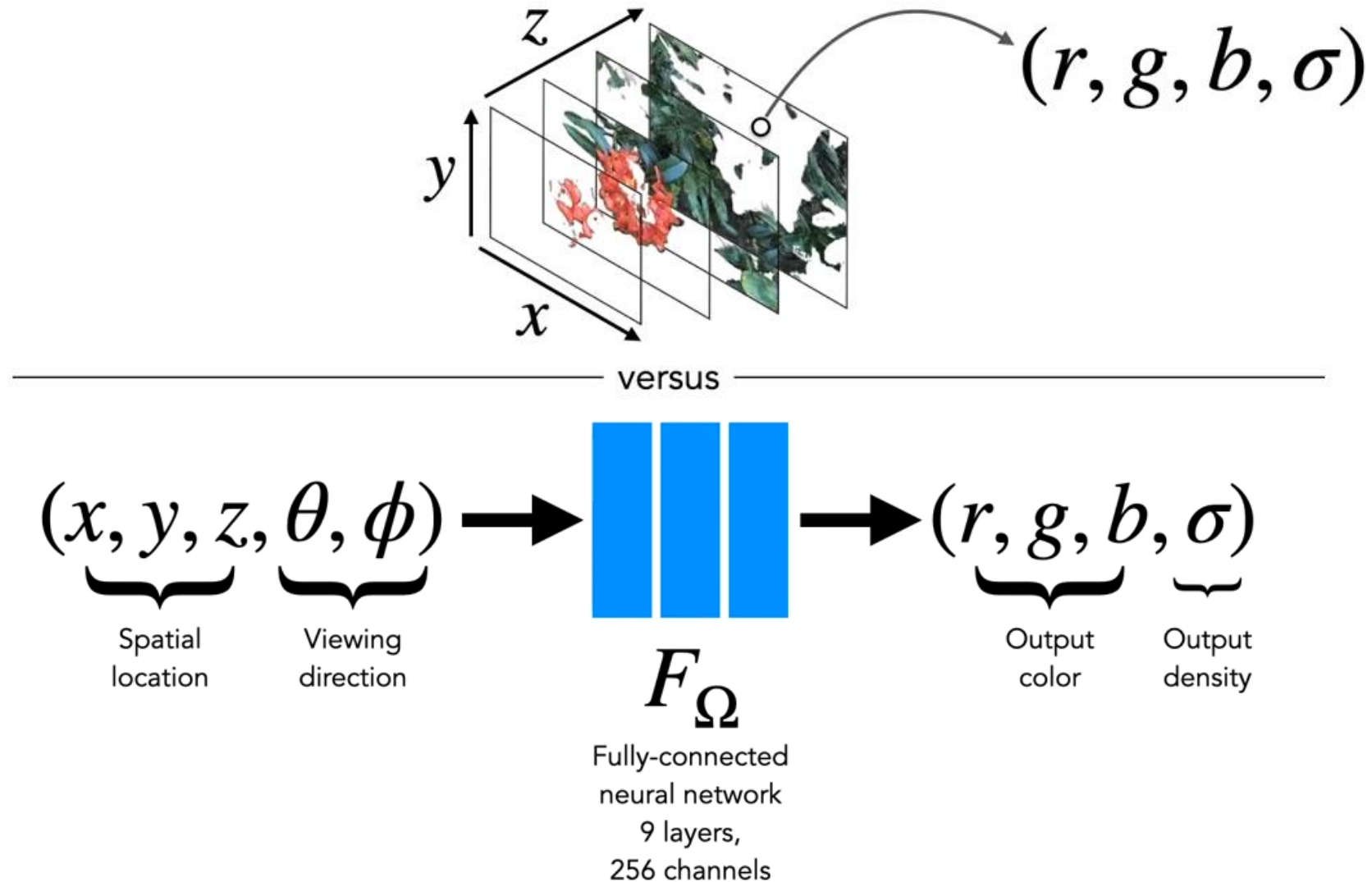
- Continuous neural network as a volumetric scene representation (5D = xyz + direction)
- Use volume rendering model to synthesize new views
- Optimize using rendering loss for one scene (no prior training)
- Apply positional encoding before passing coordinates into network to recover high frequency details

NeRF: Neural Radiance Fields

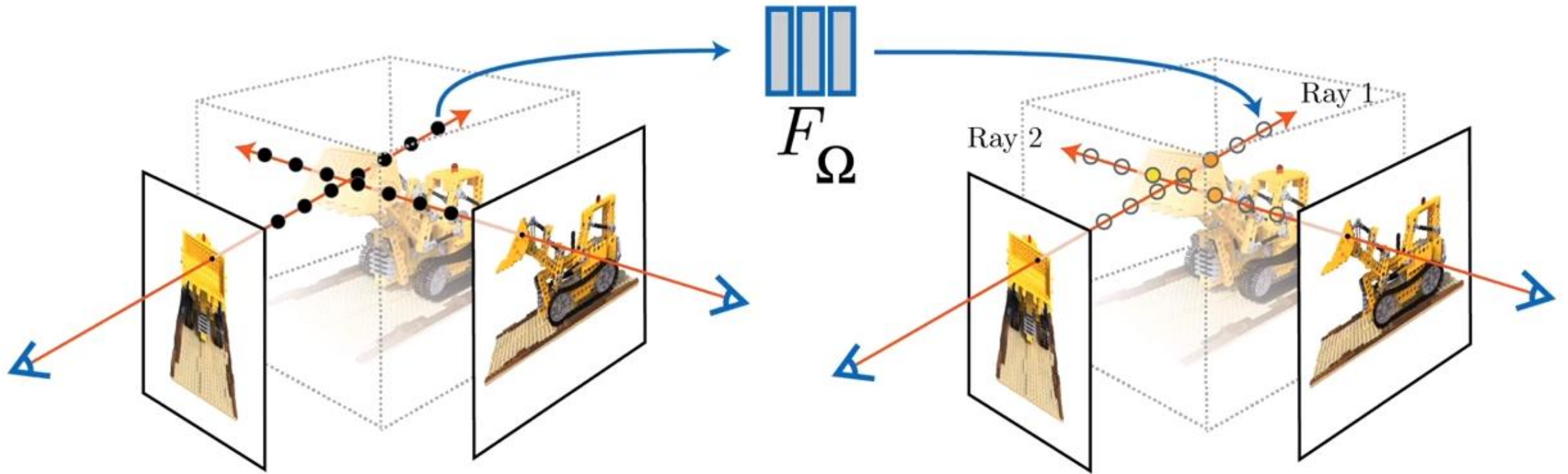


[Overview of neural radiance field scene representation and differentiable rendering procedure]

Representing a scene as a continuous 5D function



Generate views with traditional volume rendering



Generate views with traditional volume rendering

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

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

weights colors

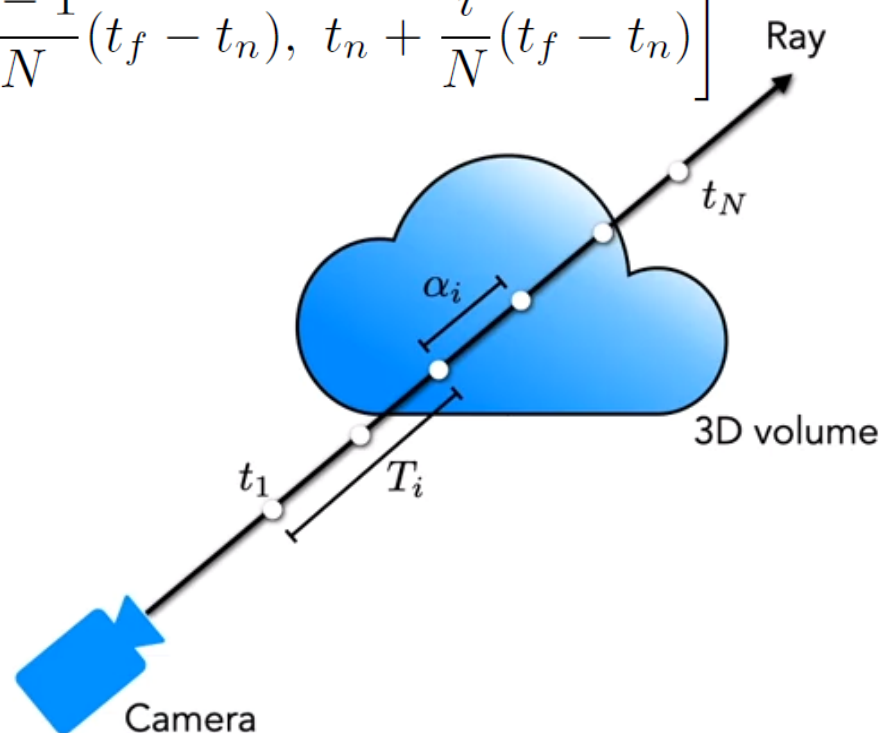
How much light is blocked earlier along ray:

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

How much light is contributed by ray segment i :

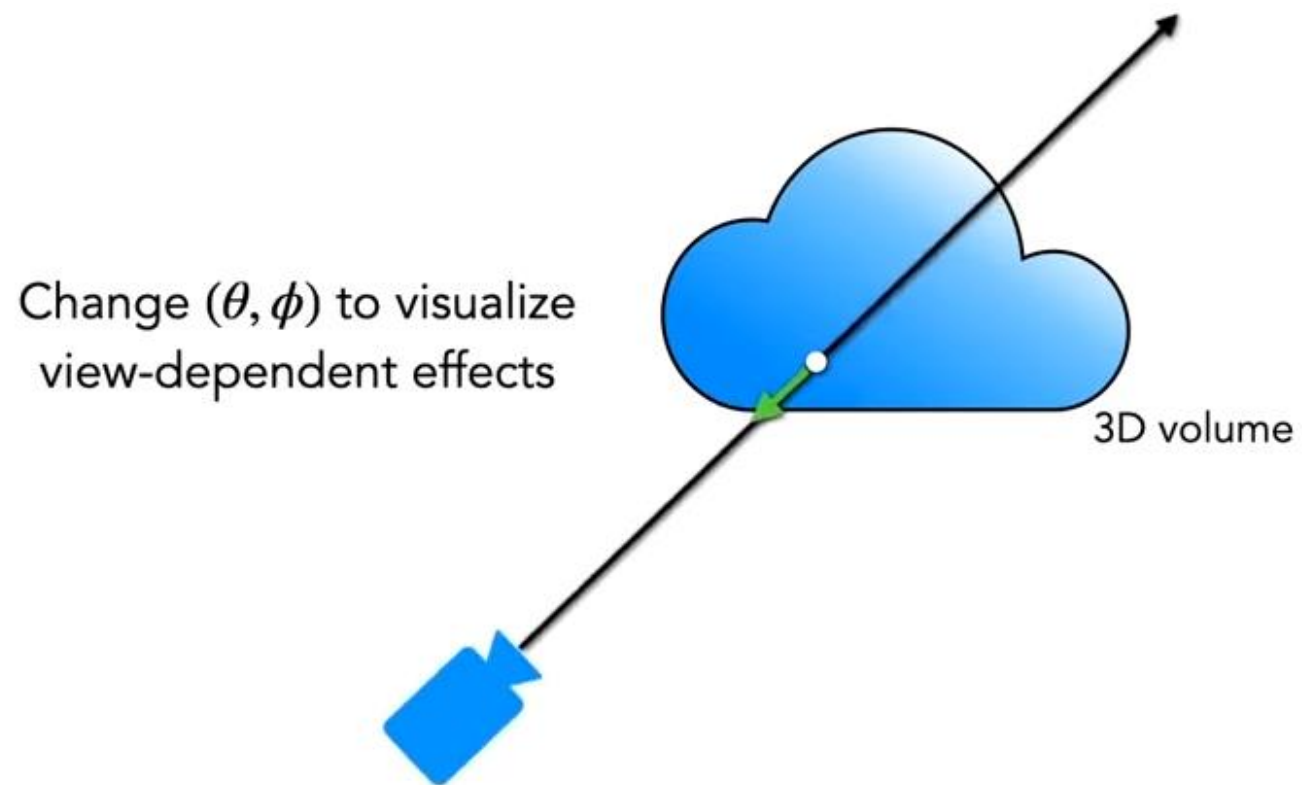
$$\alpha_i = 1 - e^{-\sigma_i \delta t_i}$$

$$t_i \sim \mathcal{U} \left[t_n + \frac{i-1}{N} (t_f - t_n), t_n + \frac{i}{N} (t_f - t_n) \right]$$



$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i, \text{ where } T_i = \exp \left(- \sum_{j=1}^{i-1} \sigma_j \delta_j \right)$$

Viewing directions as input



Viewing directions as input

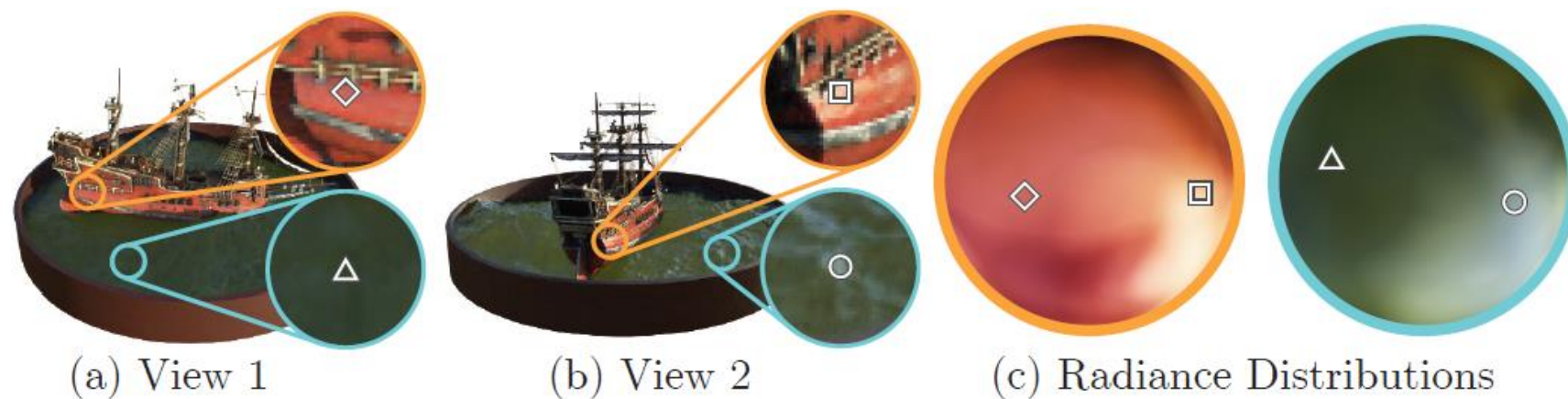
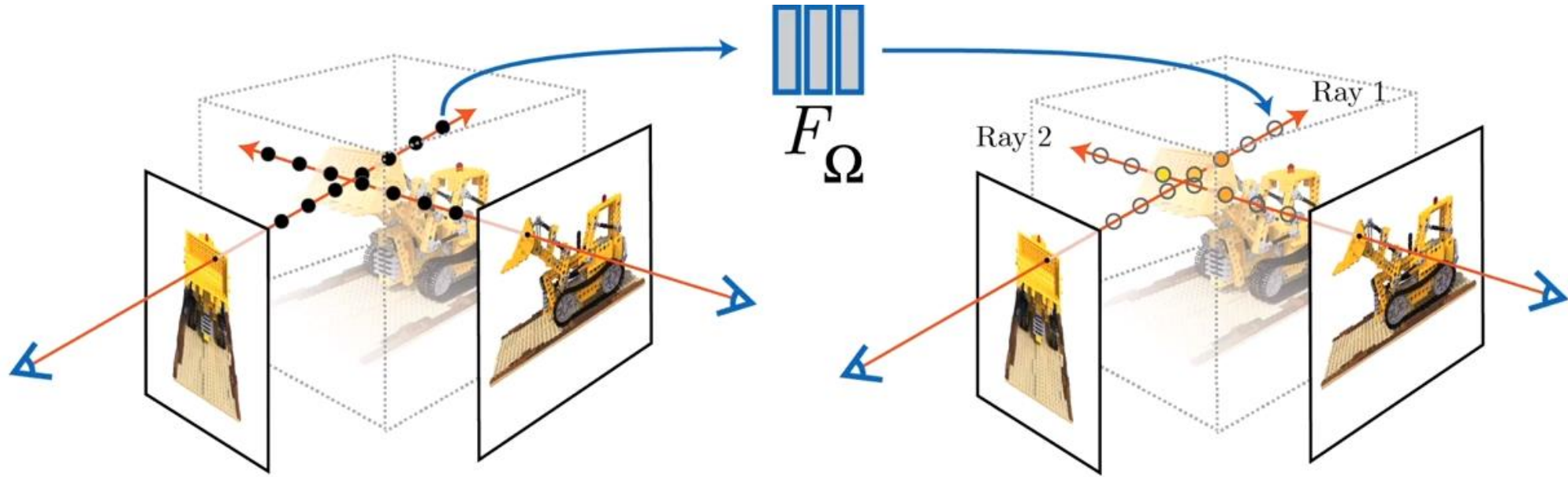


Fig. 3: A visualization of view-dependent emitted radiance. Our neural radiance field representation outputs RGB color as a 5D function of both spatial position \mathbf{x} and viewing direction \mathbf{d} . Here, we visualize example directional color distributions for two spatial locations in our neural 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.

Optimize with gradient descent on rendering loss

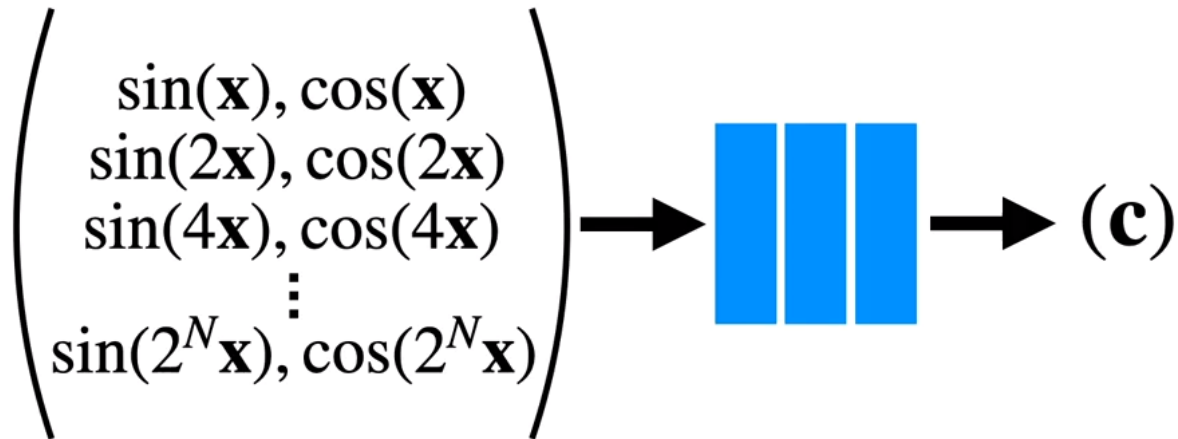


$$\min_{\Omega} \sum_i \|\text{render}^{(i)}(F_{\Omega}) - I_{\text{gt}}^{(i)}\|^2$$

Positional encoding

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

In experiments, they set $L = 10$ for $\gamma(\mathbf{x})$ and $L = 4$ for $\gamma(d)$.

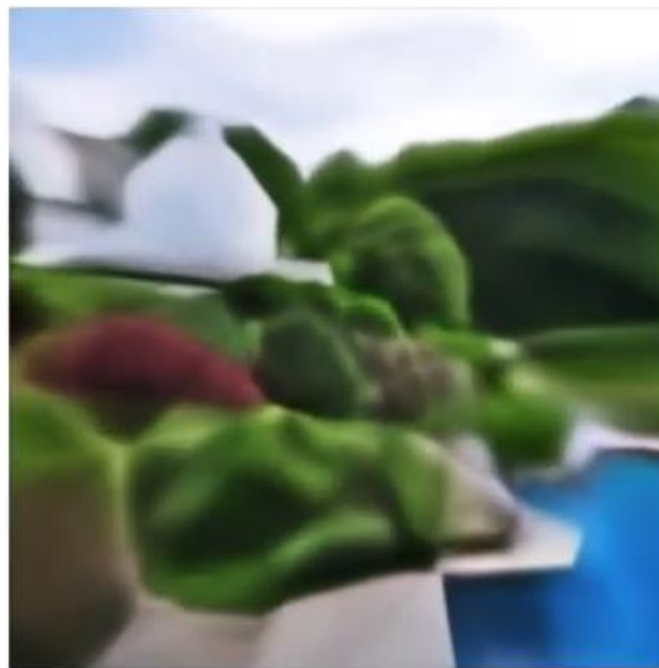


The effectiveness of positional encoding

Ground truth image



Standard fully-connected net



With "positional encoding"



Experiments: Quantitative comparison

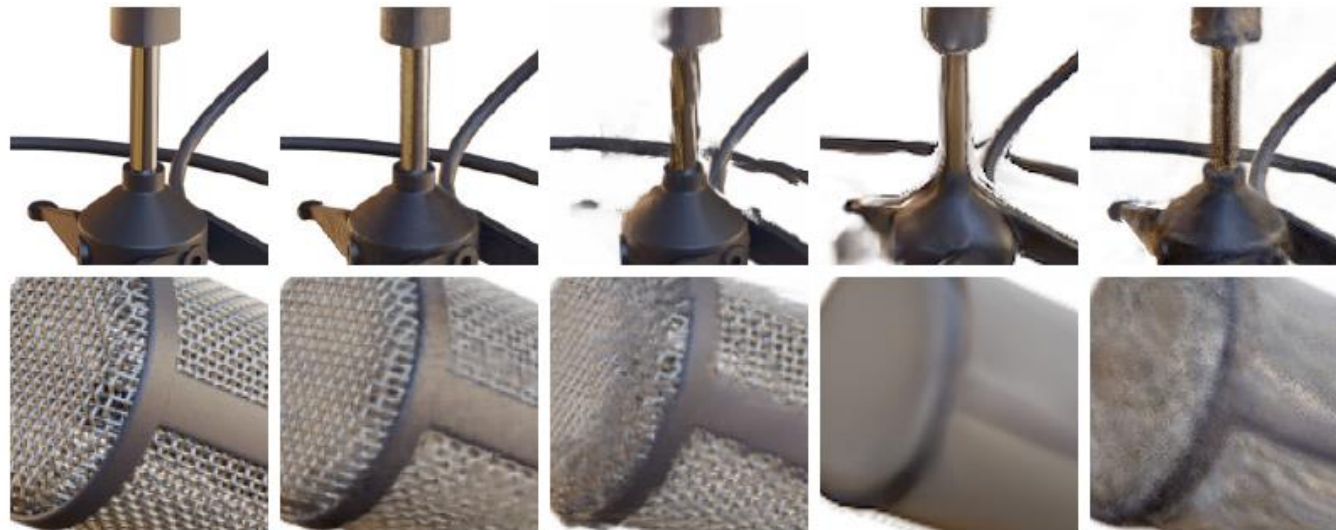
Method	Diffuse Synthetic 360° [41]			Realistic Synthetic 360°			Real Forward-Facing [28]		
	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓
SRN [42]	33.20	0.963	0.073	22.26	0.846	0.170	22.84	0.668	0.378
NV [24]	29.62	0.929	0.099	26.05	0.893	0.160	-	-	-
LLFF [28]	34.38	0.985	0.048	24.88	0.911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

Table 1: Our method quantitatively outperforms prior work on datasets of both synthetic and real images. We report PSNR/SSIM (higher is better) and LPIPS [50] (lower is better). The DeepVoxels [41] dataset consists of 4 diffuse objects with simple geometry. Our realistic synthetic dataset consists of pathtraced renderings of 8 geometrically complex objects with complex non-Lambertian materials. The real dataset consists of handheld forward-facing captures of 8 real-world scenes (NV cannot be evaluated on this data because it only reconstructs objects inside a bounded volume). Though LLFF achieves slightly better LPIPS, we urge readers to view our supplementary video where our method achieves better multiview consistency and produces fewer artifacts than all baselines.

Experiments: Qualitative comparison



Microphone



Materials

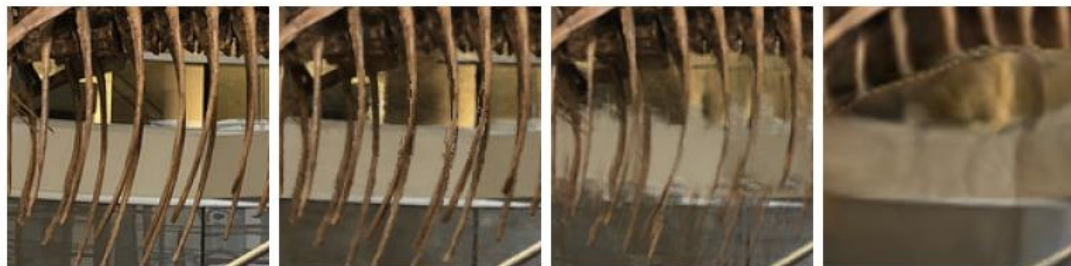


Ground Truth NeRF (ours) LLFF [28] SRN [42] NV [24]

Experiments: Qualitative comparison



T-Rex



Orchid



Ground Truth

NeRF (ours)

LLFF [28]

SRN [42]

Experiments: Ablation study

	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)	31.01	0.947	0.081

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.

Thank you!