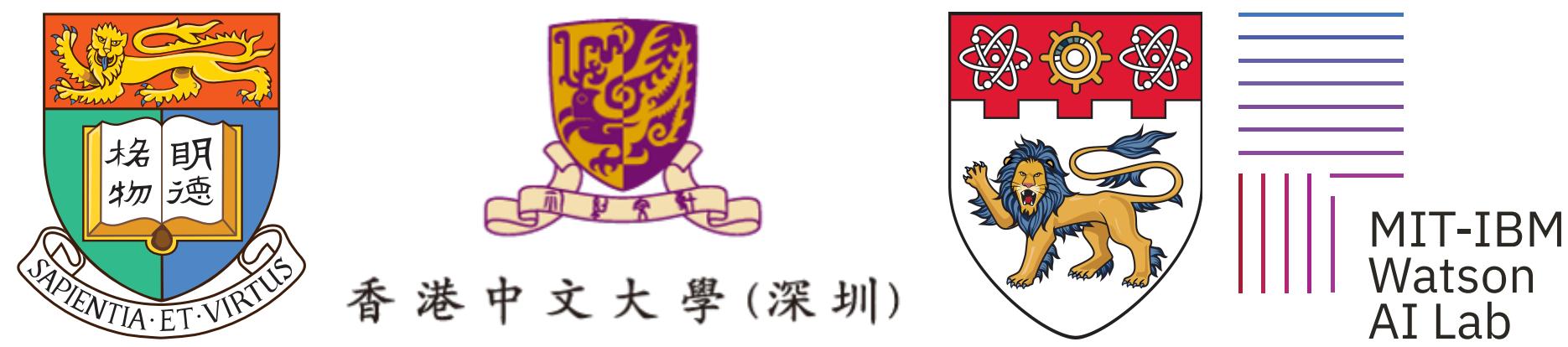


# S<sup>3</sup>-NeRF: Neural Reflectance Field from Shading and Shadow under a Single Viewpoint



Wenqi Yang<sup>1</sup> Guanying Chen<sup>2</sup> Chaofeng Chen<sup>3</sup> Zhenfang Chen<sup>4</sup> Kwan-Yee K. Wong<sup>1</sup>

<sup>1</sup>The University of Hong Kong <sup>2</sup>FNii and SSE, CUHK-Shenzhen

<sup>3</sup>Nanyang Technological University <sup>4</sup>MIT-IBM Watson AI Lab



Code & Data & Model:

<https://yqg.github.io/s3nerf/>



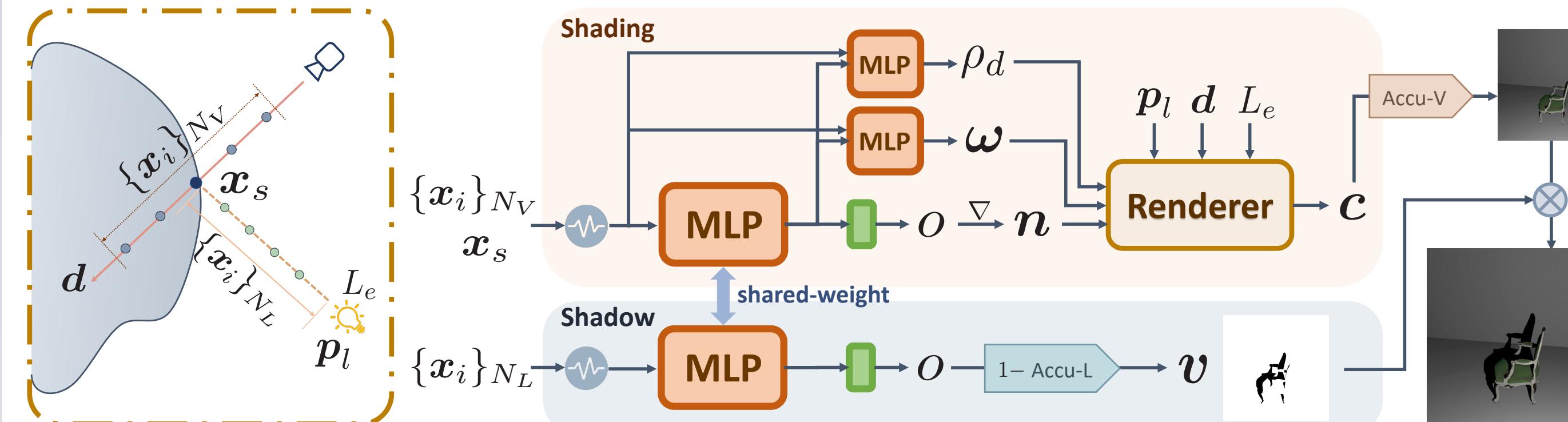
## Idea

- “Dual problem” of Multi-View Scene reconstruction – utilize **single-view** images captured under **different point lights** to reconstruct a scene.
- Existing single-view methods – only recover a **2.5D** scene representation (*i.e.*, a normal / depth map for the visible surface).
- Ours – Learn a **3D** neural reflectance field.
- MVS – Rely on multi-view photo-consistency to infer scene geometry.
- Ours – Exploit two information-rich monocular cues – **shading & shadow**.

## Method

### Overview

Given  $N$  images captured from a single viewpoint under different near point lights, S<sup>3</sup>-NeRF targets at recovering the geometry and materials of the scene.



We first apply root-finding to locate the surface intersection point  $x_s$ .

- $N_V$  points on the camera ray are sampled around the surface to generate accumulated shading values.
- $N_L$  points are sampled on the surface-to-light segment to calculate the light visibility ( $O(N_L)$  MLP queries per ray, calculated in an online manner).

### Rendering Model:

We consider non-Lambertian surfaces with spatially-varying BRDFs. The rendering equation for a surface point  $\mathbf{x}$  viewed from a direction  $\mathbf{d}$  under a near point light  $(\mathbf{p}_l, L_e)$  can be written as

$$f_c(\mathbf{d}, \mathbf{p}_l, L_e; \mathbf{x}) = \underbrace{L_{int}(\mathbf{p}_l, L_e; \mathbf{x})}_{\text{Light Intensity}} \underbrace{f_m(\mathbf{d}, \mathbf{w}_i(\mathbf{p}_l; \mathbf{x}); \mathbf{x})}_{\text{BRDF Value}} \max(\mathbf{w}_i(\mathbf{p}_l; \mathbf{x}) \cdot \mathbf{n}(\mathbf{x}), 0) \cdot \underbrace{\text{Shading}}_{\text{Shading}}$$

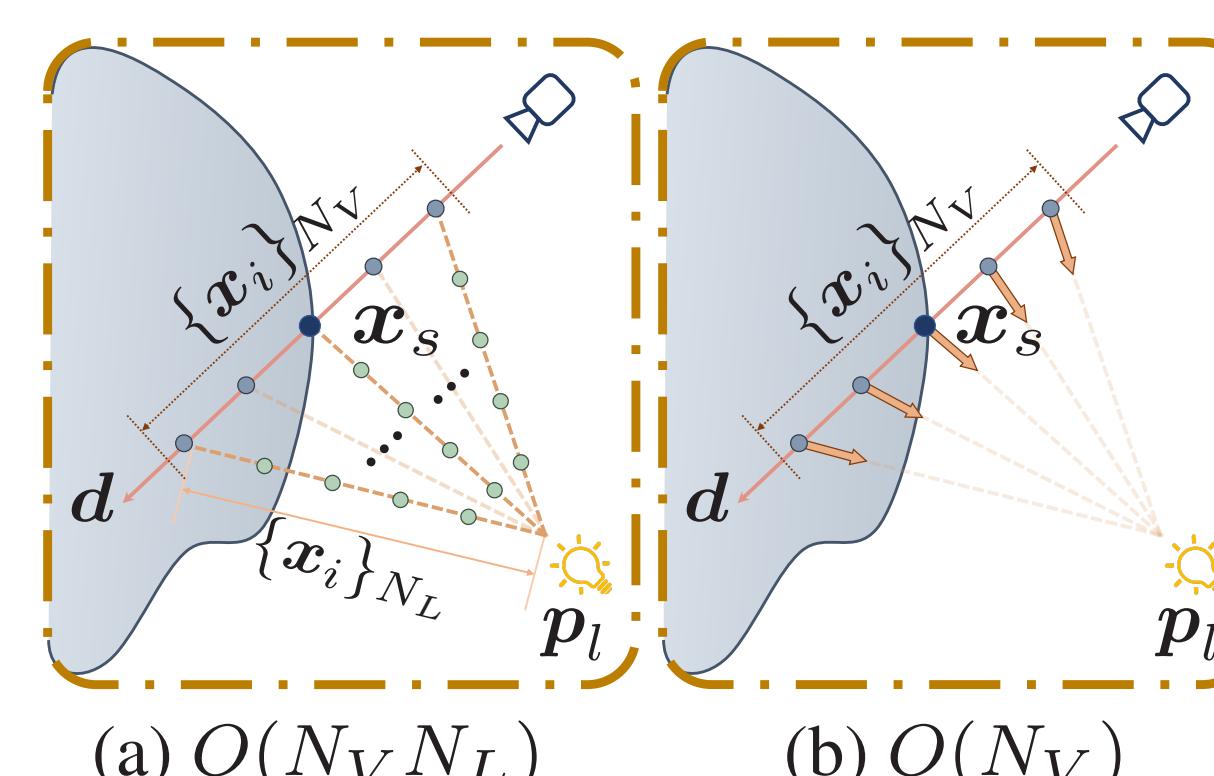
We propose to adopt a joint volume and surface rendering strategy:

$$C_v(\mathbf{r}) = f_v(\mathbf{p}_l; \mathbf{x}_s) \sum_{i=1}^{N_V} o(\mathbf{x}_i) \prod_{j < i} (1 - o(\mathbf{x}_j)) f_c(\mathbf{x}_i, \mathbf{d}, \mathbf{p}_l, L_e),$$

$$C_s(\mathbf{r}) = f_v(\mathbf{p}_l; \mathbf{x}_s) f_c(\mathbf{d}, \mathbf{p}_l, L_e; \mathbf{x}_s).$$

### Alternative Shadow Modeling:

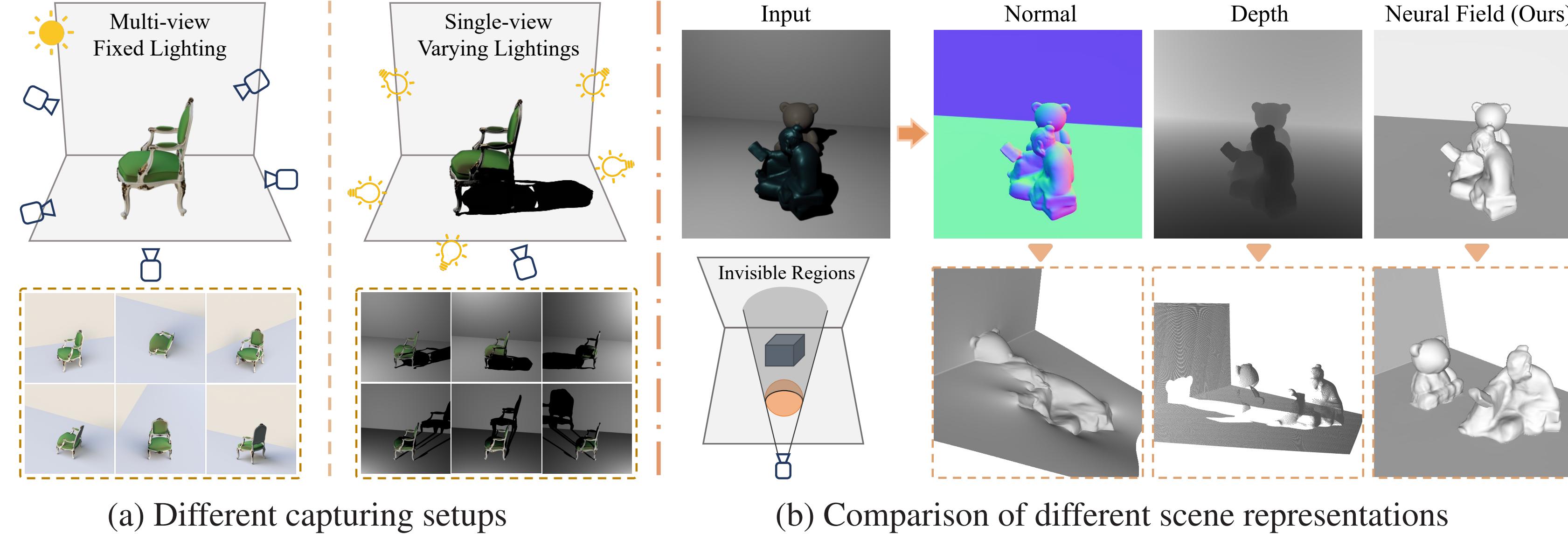
- Calculate light visibilities for all  $N_V$  points sampled along the ray (computationally expensive).
- Adopt an MLP to directly regress light visibility of a point to reduce the queries for each ray.



## Overview

### Contributions:

- Novel problem: 3D neural reflectance field optimization from single-view images captured under different point lights.
- Exploit monocular shading and shadow cues to jointly recover the geometry and BRDFs of a scene.
- Adopt an efficient online shadow computation to fully exploit the information-rich shading and shadow cues.
- Experiments on multiple challenging datasets show that S<sup>3</sup>-NeRF can faithfully reconstruct a complete scene geometry from single-view images and is robust to depth discontinuity.



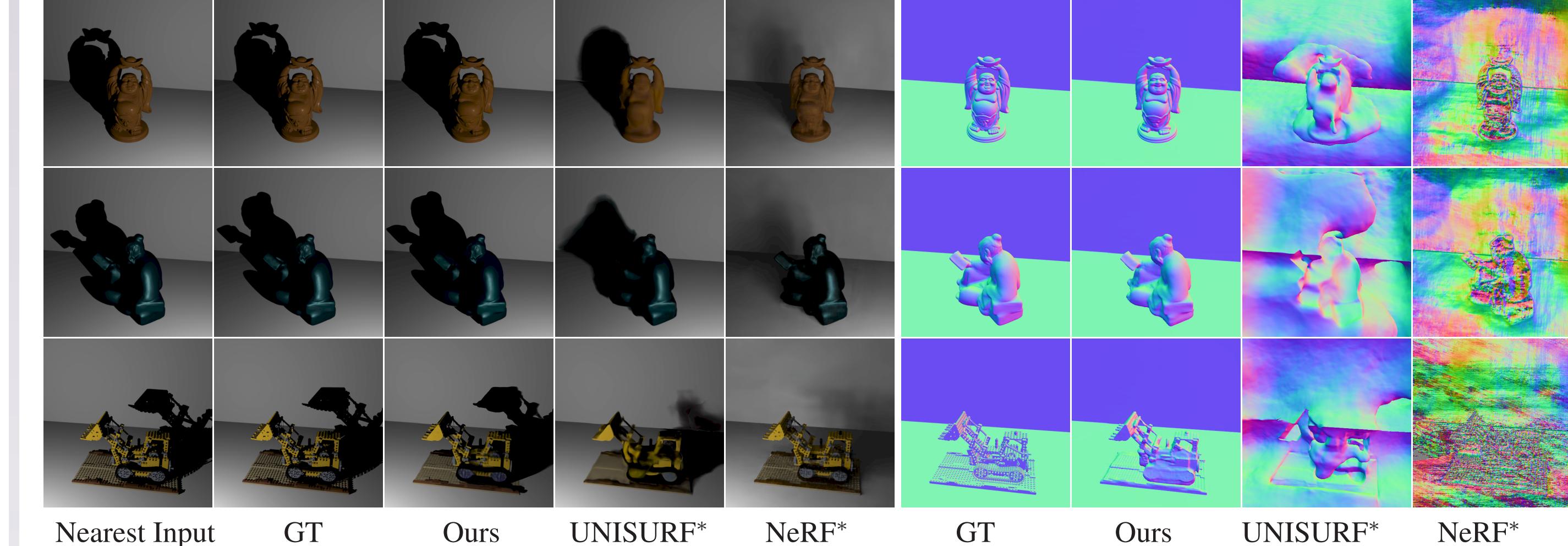
## Experiments & Results

### Comparison with Neural Radiance Field Methods:

#### Quantitative results on relighting and normal estimation

Method	BUDDHA		READING		BUNNY		CHAIR		LEGO		HOTDOG	
	PSNR↑	MAE↓										
NeRF* [1]	38.57	70.12	39.50	72.60	37.41	68.35	35.25	88.46	35.56	91.09	39.80	72.07
UNISURF* [2]	41.51	54.86	40.54	60.59	38.48	54.27	34.98	47.79	34.55	45.81	38.64	51.00
Ours	<b>43.42</b>	<b>2.44</b>	<b>43.13</b>	<b>2.03</b>	<b>40.43</b>	<b>1.72</b>	<b>36.33</b>	<b>1.83</b>	<b>35.54</b>	<b>6.49</b>	<b>38.01</b>	<b>2.50</b>

#### Qualitative results on relighting and normal estimation

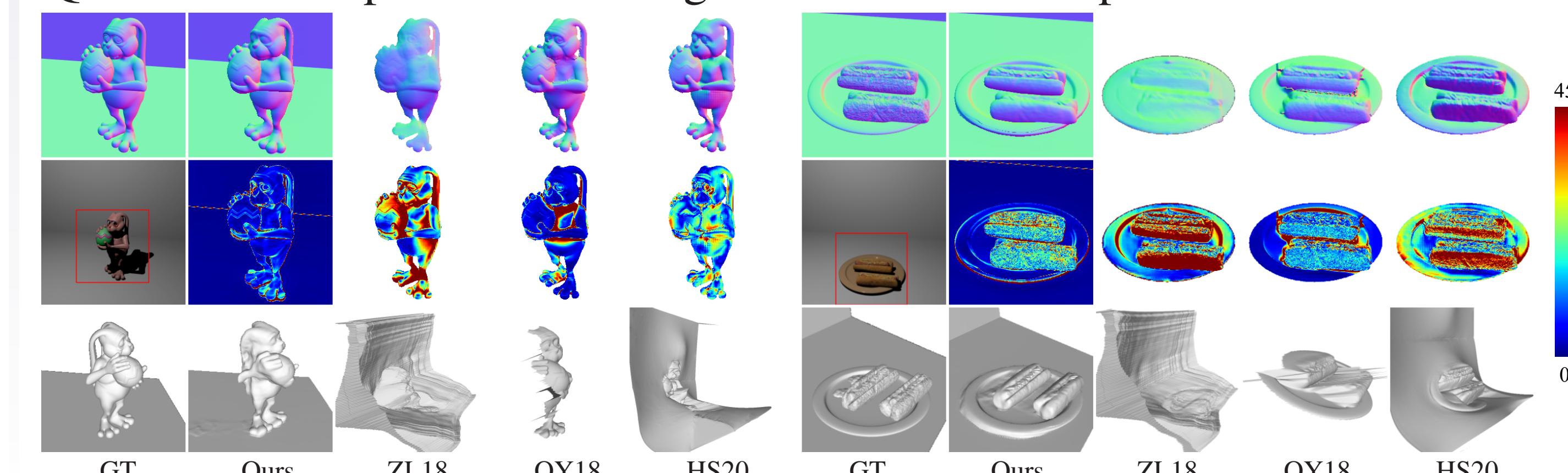


### Comparison with Single-view Shape Estimation Methods:

#### Quantitative comparison (only object regions)

Method	BUDDHA		READING		BUNNY		CHAIR		LEGO		HOTDOG	
	MAE↓	Depth L1↓	MAE↓	Depth L1↓	MAE↓	Depth L1↓	MAE↓	Depth L1↓	MAE↓	Depth L1↓	MAE↓	Depth L1↓
ZL18 [3]	37.51	19.84	37.29	25.97	31.40	17.68	39.53	41.19	46.82	34.56	39.74	18.02
QY18 [4]	<b>12.25</b>	3.81	40.84	26.13	14.21	4.10	29.68	15.95	33.08	17.87	16.81	8.98
HS20 [5]	18.39	6.47	27.11	18.94	16.92	10.96	29.56	13.99	33.54	13.27	27.25	13.22
Ours	14.24	<b>1.50</b>	<b>7.00</b>	<b>2.09</b>	<b>9.40</b>	<b>1.63</b>	<b>17.43</b>	<b>4.74</b>	<b>31.13</b>	<b>7.31</b>	<b>14.65</b>	<b>1.68</b>

#### Qualitative comparison with single-view normal / depth estimation baselines

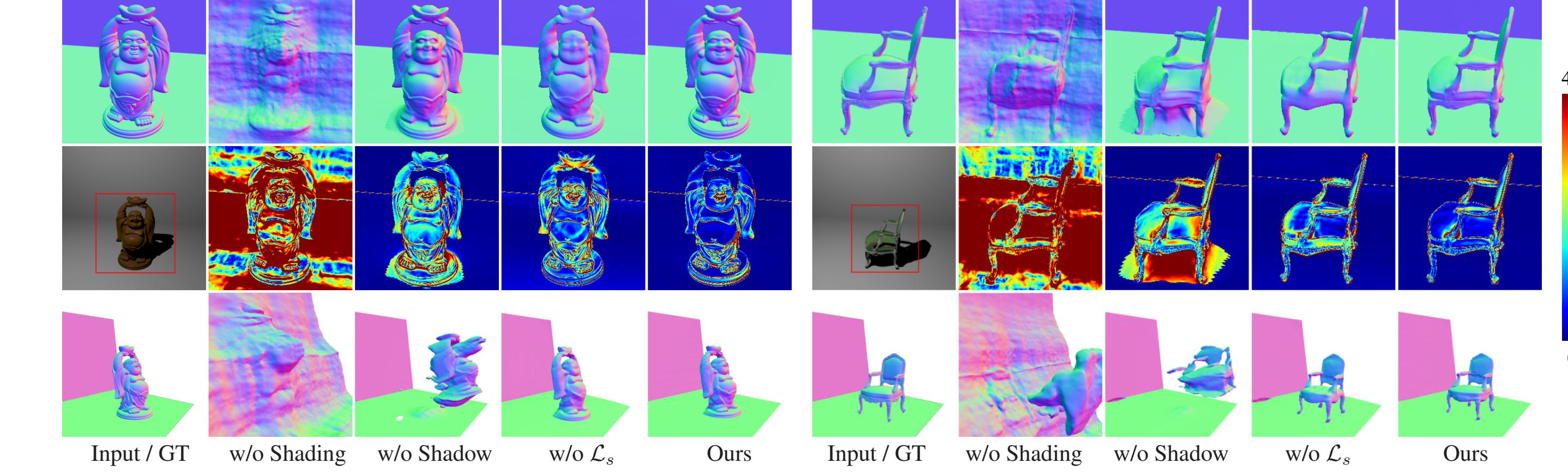


### Ablation Study:

#### Quantitative results for the ablation study

Method	CHAIR		BUNNY		BUDDHA		CHAIR		BUNNY		BUDDHA	
	MAE↓	PSNR↑										
w/o shading	32.49	33.71	40.18	38.72	35.68	41.43	—	—	—	—	—	—
w/o shadow	3.39	30.81	2.26	33.45	3.33	34.30	12.24	22.31	11.93	24.43	16.60	23.27
w/o $\mathcal{L}_s$	2.48	35.85	2.75	39.73	3.77	43.04	<b>5.10</b>	<b>28.58</b>	6.27	29.11	8.50	28.61
Ours	<b>1.83</b>	<b>36.33</b>	<b>1.72</b>	<b>40.43</b>	<b>2.44</b>	<b>43.42</b>	5.45	26.82	<b>6.11</b>	<b>29.55</b>	<b>6.89</b>	<b>31.53</b>

#### Visual results for the ablation study



### Results on Real Scenes:



### Capturing Setup:



### References:

- [1] NeRF [Mildenhall et al., ECCV20]
- [2] UNISURF [Oechsle et al., ICCV21]
- [3] ZL18 [Li et al., TOG18]
- [4] QY18 [Quéau et al., JMIV18]
- [5] HS20 [Santo et al., ECCV20]