

# LidaRF: Delving into Lidar for Neural Radiance Field on Street Scenes

Shanlin Sun<sup>1</sup>, Bingbing Zhuang<sup>3</sup>, Ziyu Jiang<sup>3</sup>, Buyu Liu<sup>3</sup>, Xiaohui Xie<sup>1</sup>, Manmohan Chandraker<sup>2,3</sup>

<sup>1</sup>University of California, Irvine   <sup>2</sup>University of California, San Diego   <sup>3</sup>NEC Labs America



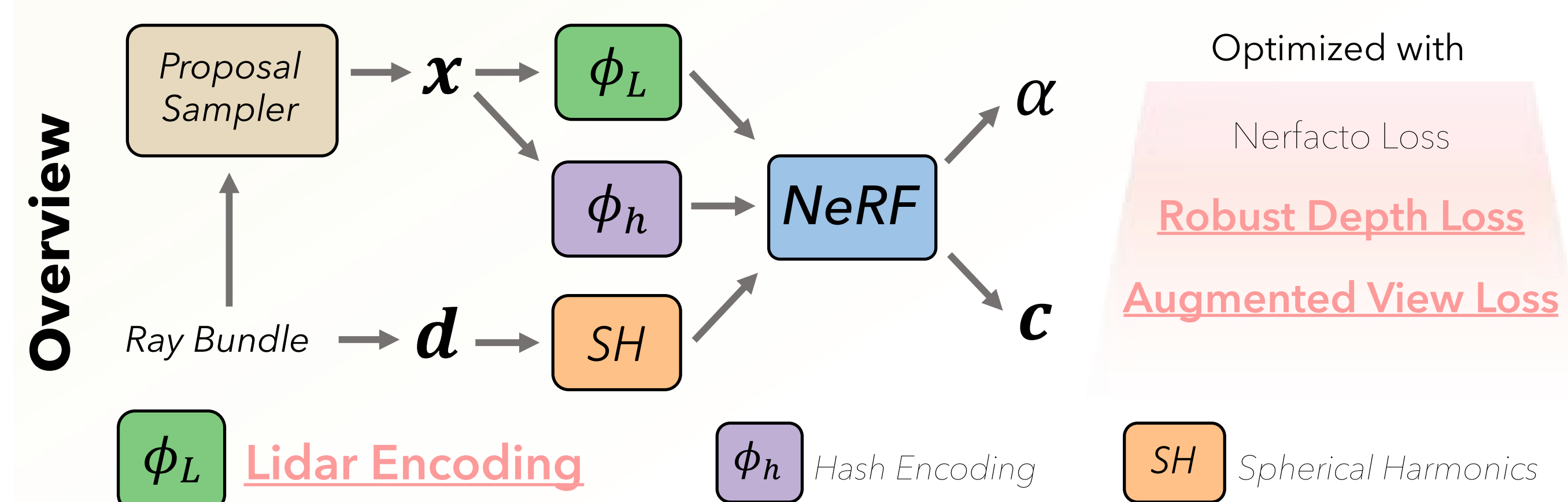
## Introduction: Appearance Simulation

Input: Front RGB camera, Lidar, sensor poses

Goal: Photorealistic appearance simulation of street scenes

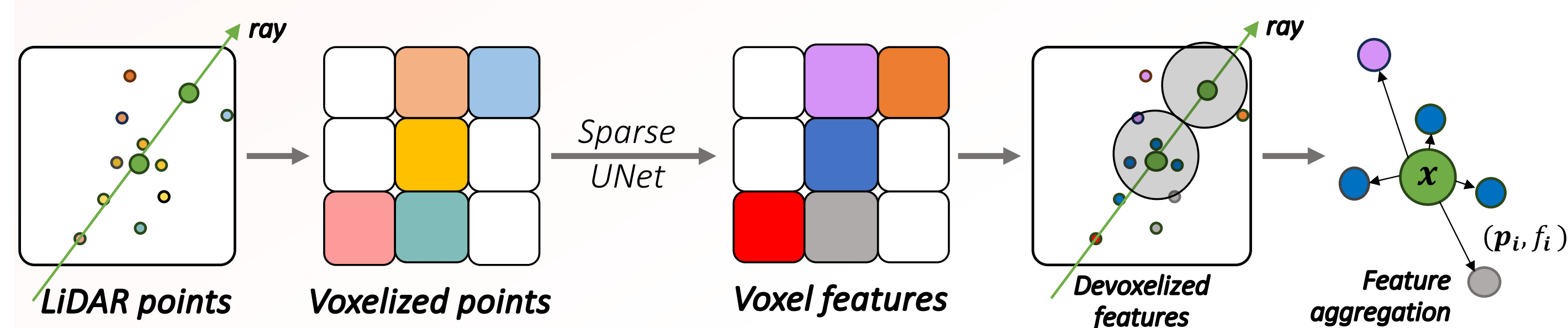
Challenge: Sparse view coverage, low-texture road surface ...

Motivation: Modern autonomous systems are often equipped with Lidar. How can we use it more than just as a depth loss?



## Contrib #1: Lidar Encoding

- Lidar holds strong potential for geometric guidance
- Lidar encoding through 3D sparse CNN has proven powerful in 3D detection, but is underexplored in NeRF
- Fuse Lidar encoding and hash grid feature



## Contrib #2: Robust Depth Supervision

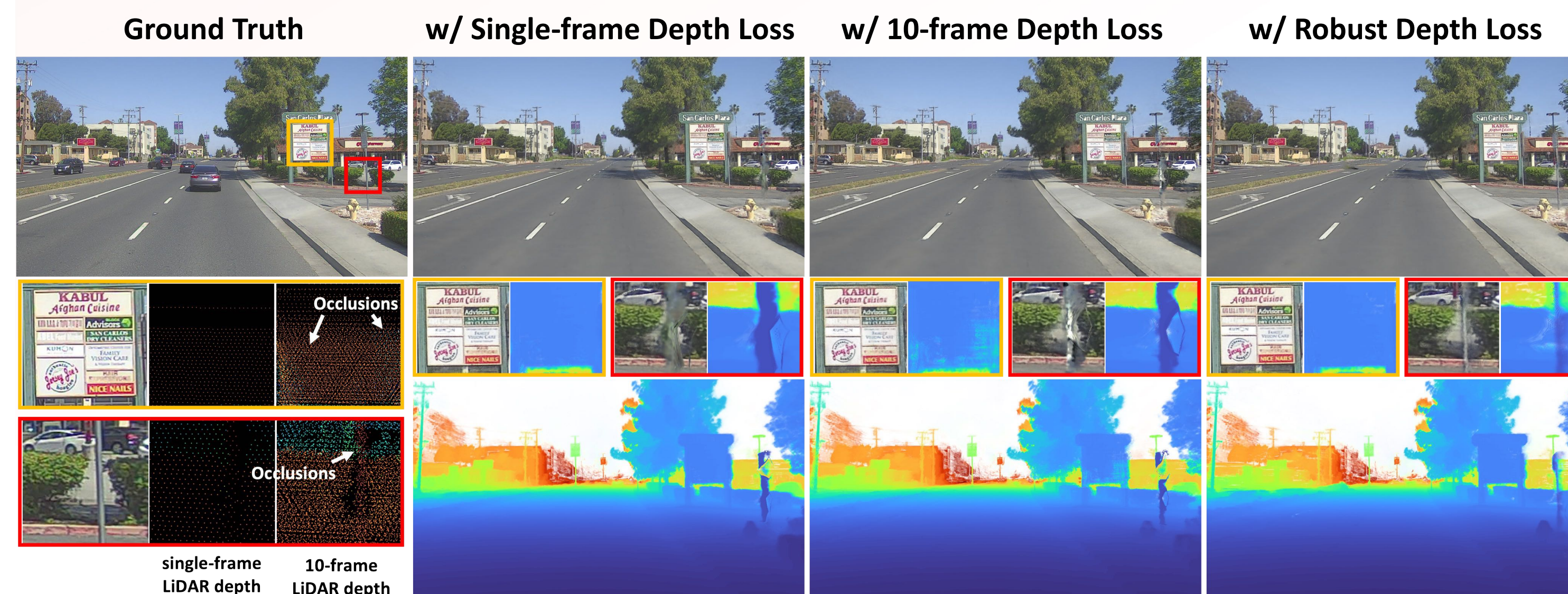
- Accumulate adjacent Lidar frames for denser supervision
- Inter-points occlusion due to the camera-Lidar displacement
- For samples in  $\mathcal{D}_{\text{reliable}}^m$ , we adopt the pixel-level depth loss

$$\mathcal{D}_{\text{reliable}}^m = \{ \mathcal{D}_i \mid \mathcal{D}_i \leq \epsilon_t^m, \mathcal{D}_i \leq \hat{\mathcal{D}}_i + \epsilon_o^m, \mathcal{D}_i \in \mathcal{D} \}$$

$$\epsilon_t^m = \min\{\alpha_t \epsilon_t^{m-1}, \epsilon_t\}, \quad \alpha_t > 1$$

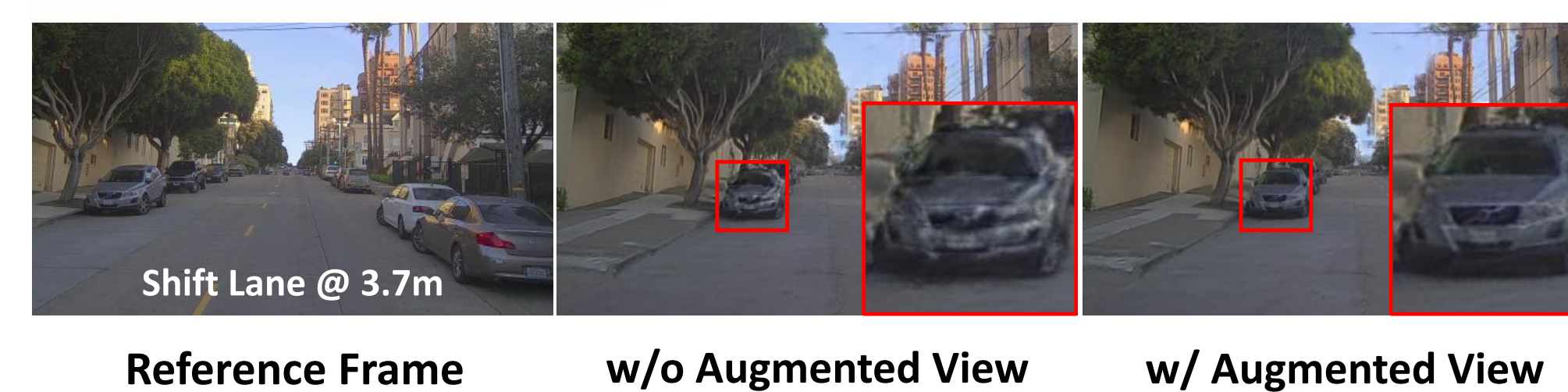
$$\epsilon_o^m = \max\{\alpha_o \epsilon_o^{m-1}, \epsilon_o\}, \quad \alpha_o < 1$$

Intuition: count on near points initially and gradually add far points congruent along with NeRF training



## Contrib #3: Augmented View Supervision

- Project Lidar points to "render" more training views
- Handle occlusions using the same scheme as "contribution #2"



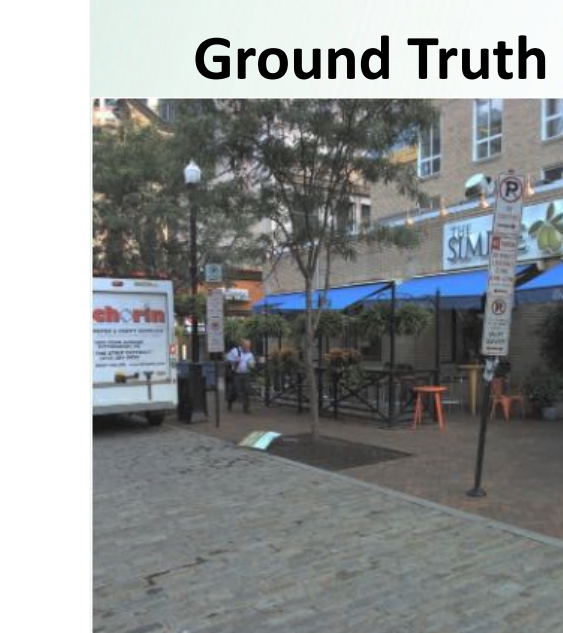
Future work: in-paint sparse projection with diffusion priors

## Experiment: Comparisons with SOTA

PandaSet Dataset:

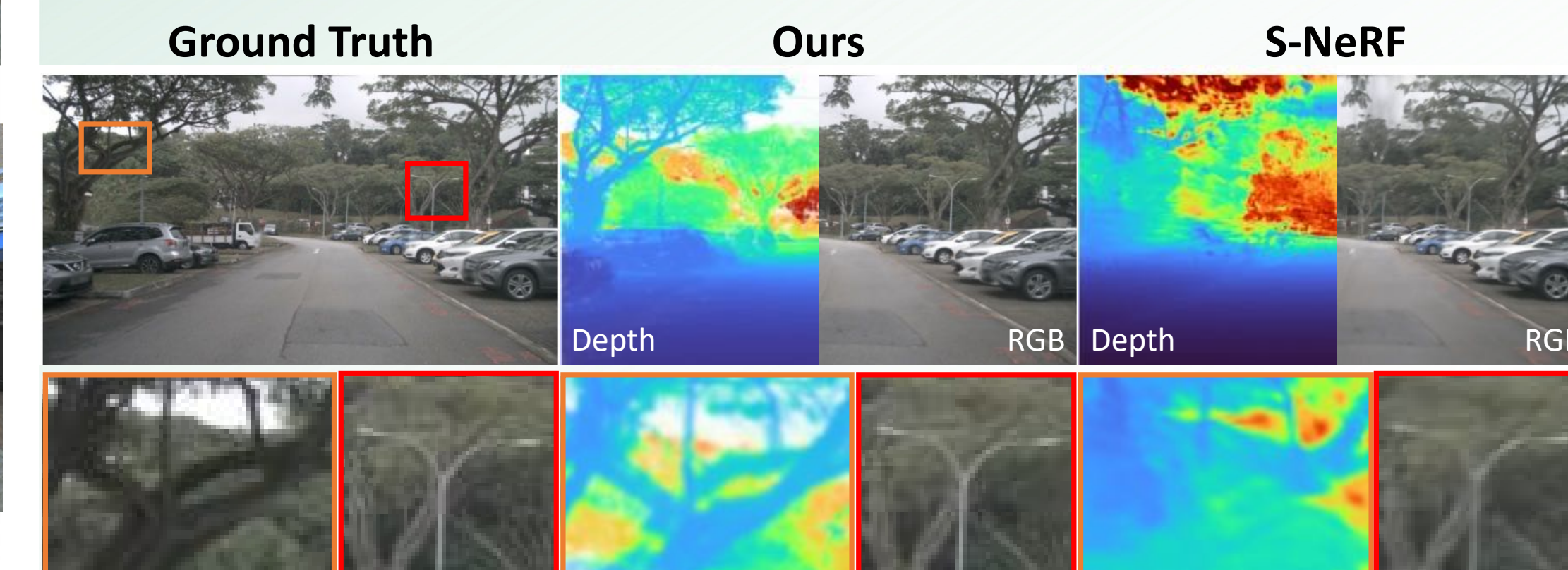


Argoverse Dataset:



Methods	Interpolation			Lane Shift	
	PSNR↑	SSIM↑	LPIPS↓	FID↓ @ 2m	FID↓ @ 3.7m
Instant-NGP	24.282	0.733	0.408	140.3	173.2
Mip-NeRF 360	23.693	0.691	0.496	189.4	231.1
Nerfatto	27.122	0.804	0.268	116.7	151.0
UniSim	26.014	0.768	0.342	118.5	141.3
Ours	27.255	<b>0.812</b>	<b>0.224</b>	<b>106.5</b>	<b>126.0</b>

NuScenes Dataset:



Methods	S-NeRF	Ours			
		w/o $\mathcal{L}_{ds}$	w/o $\phi_L$	w/o $\mathcal{L}_{aug}$	Full
PSNR↑	29.377	30.629	31.001	31.133	<b>31.162</b>
SSIM↑	0.859	0.871	0.873	0.883	<b>0.884</b>
LPIPS↓	0.349	0.278	0.237	0.222	<b>0.211</b>

