



清华大学电子工程系

Department of Electronic Engineering, Tsinghua University



美团



# Rad-NeRF: Ray-decoupled Training of Neural Radiance Field

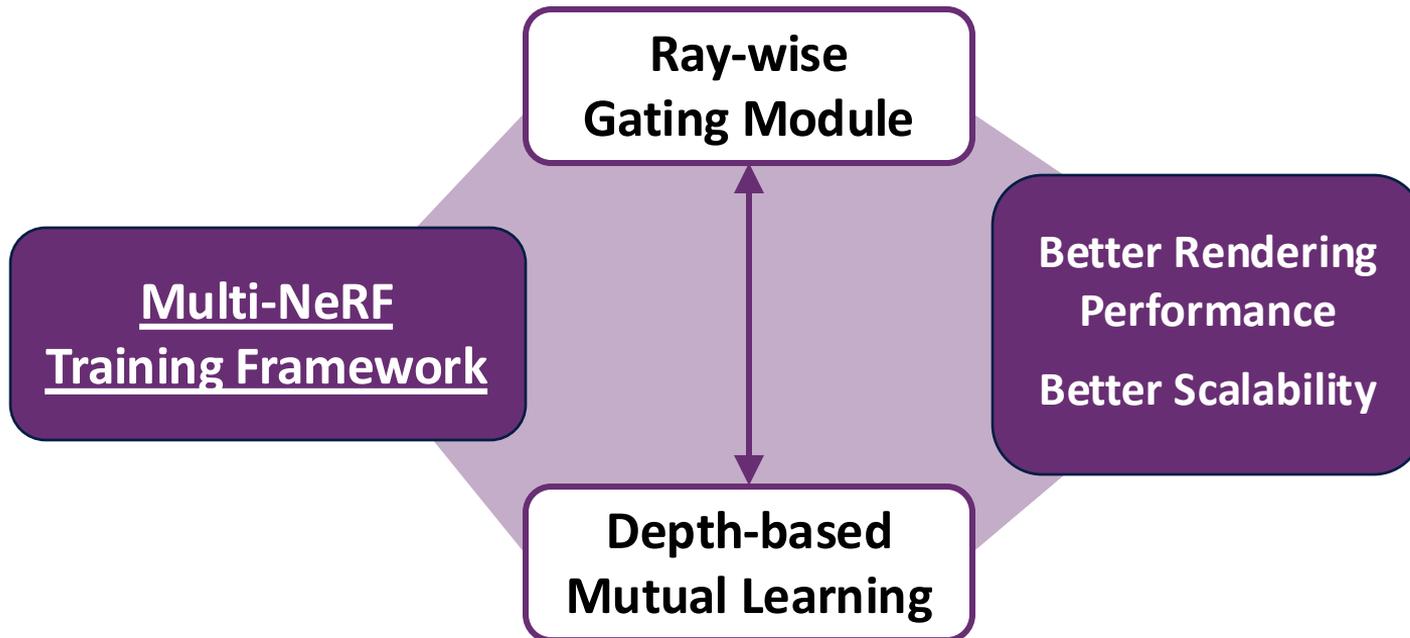
Lidong Guo<sup>1\*</sup>, Xuefei Ning<sup>1\*</sup>, Yonggan Fu<sup>2</sup>, Tianchen Zhao<sup>3</sup>, Zhuoliang Kang<sup>3</sup>,  
Jincheng Yu<sup>1</sup>, Yingyan (Celine) Lin<sup>2</sup>, Yu Wang<sup>1</sup>

<sup>1</sup>Tsinghua University, <sup>2</sup>Georgia Institute of Technology, <sup>3</sup>Meituan

E-mail: [gld21@mails.tsinghua.edu.cn](mailto:gld21@mails.tsinghua.edu.cn), [foxdoraame@gmail.com](mailto:foxdoraame@gmail.com), [yu-wang@tsinghua.edu.cn](mailto:yu-wang@tsinghua.edu.cn)



## Breaking the limitation of NeRF Scalability in Complex Scenes

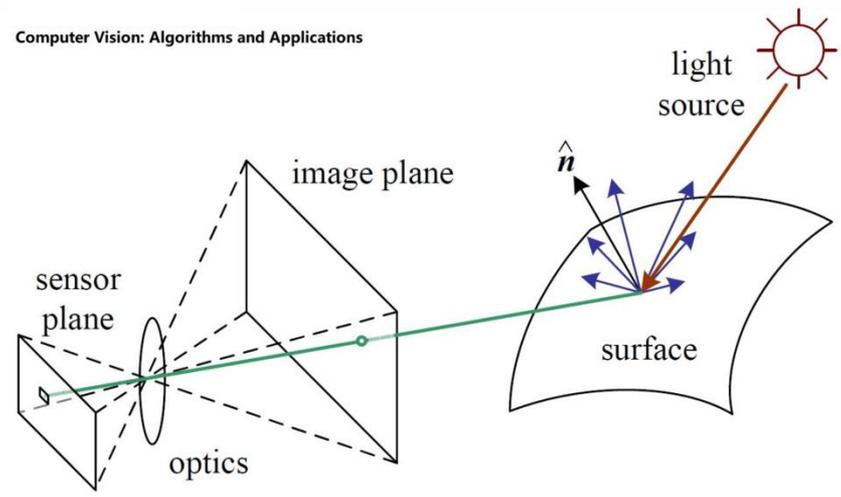




## ➤ Render:

Transform a scene representation (camera, light, surface geometry, etc.) into static images

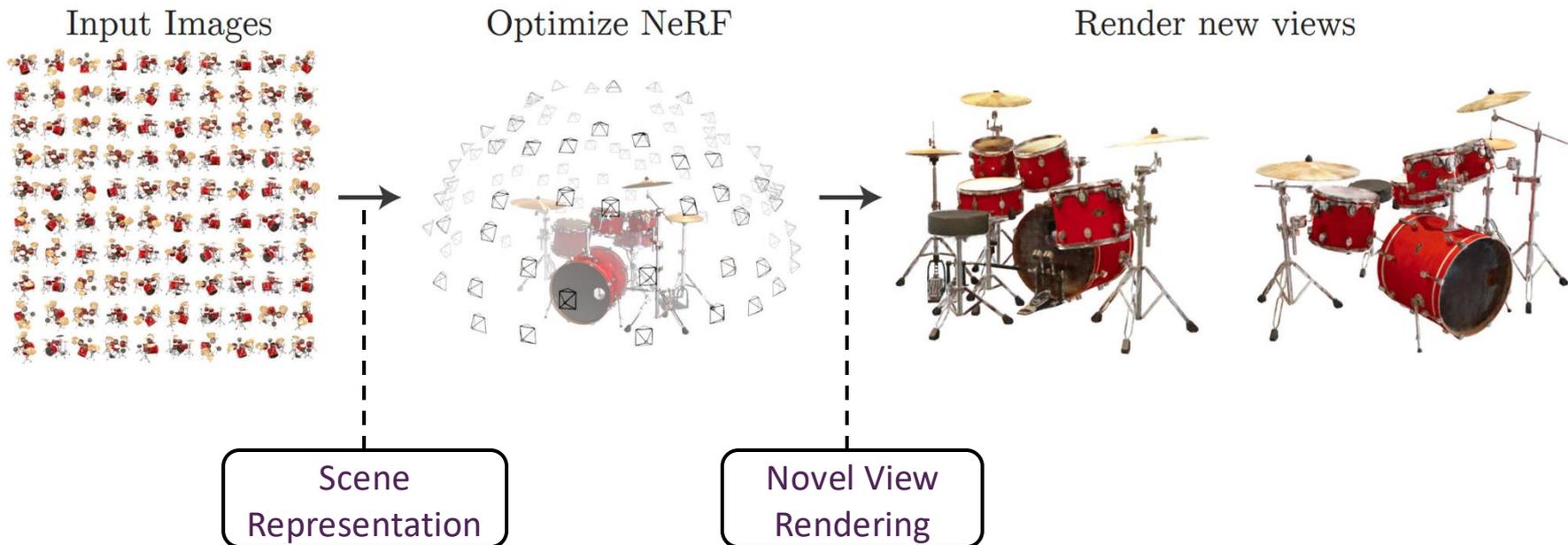
*“Taking pictures with a computer-simulated camera, provided that a 3D representation of the scene already exists”*



The physical process of taking pictures

## ➤ Rendering Task: Novel View Synthesis

Given the source image and source pose, as well as the unseen target pose, render and generate the image corresponding to the target pose.



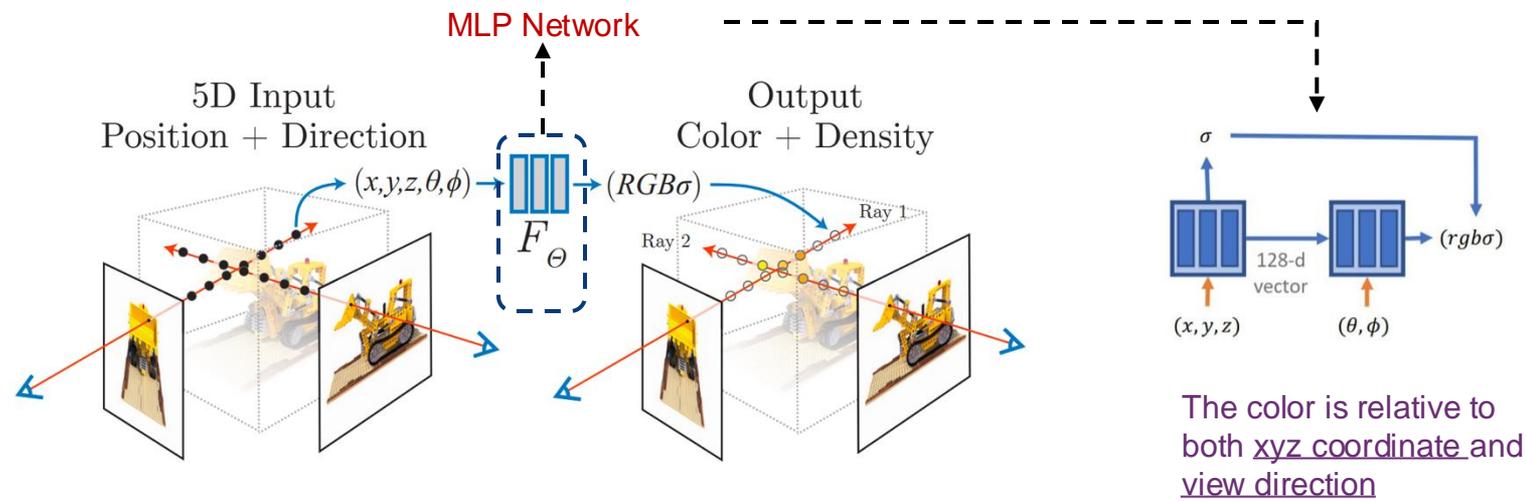


# Neural Radiance Field



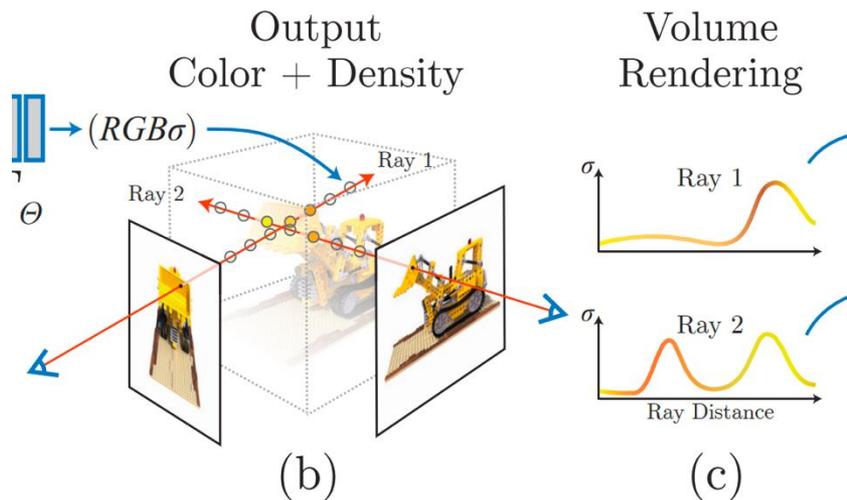
## ➤ Neural Network is adopted to fit scene representation

- A ray is made up of infinite number of 3D points in target scene
- The NN encodes the 3-d coordination and 2-d direction, outputs the color/density of each point



## ➤ Neural Network is adopted to fit scene representation

- A ray is made up of infinite number of 3D points in target scene
- The NN encodes the 3-d coordination and 2-d direction, outputs the color/density of each point
- Integrate the information of each point on the ray to obtain Pixel (Ray) color (**Volume Rendering**)



$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))\mathbf{c}(\mathbf{r}(t), \mathbf{d})dt,$$

$$T(t) = \exp\left(-\int_{t_n}^t \sigma(\mathbf{r}(s))ds\right) \text{ Transmittance}$$

Convert the integral process to discrete summation

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N T_i(1 - \exp(-\sigma_i\delta_i))\mathbf{c}_i, \quad \text{Sampling Interval}$$

$$T_i = \exp\left(-\sum_{j=1}^{i-1} \sigma_j\delta_j\right)$$

NeRF still exhibits rendering defects on complex scenes



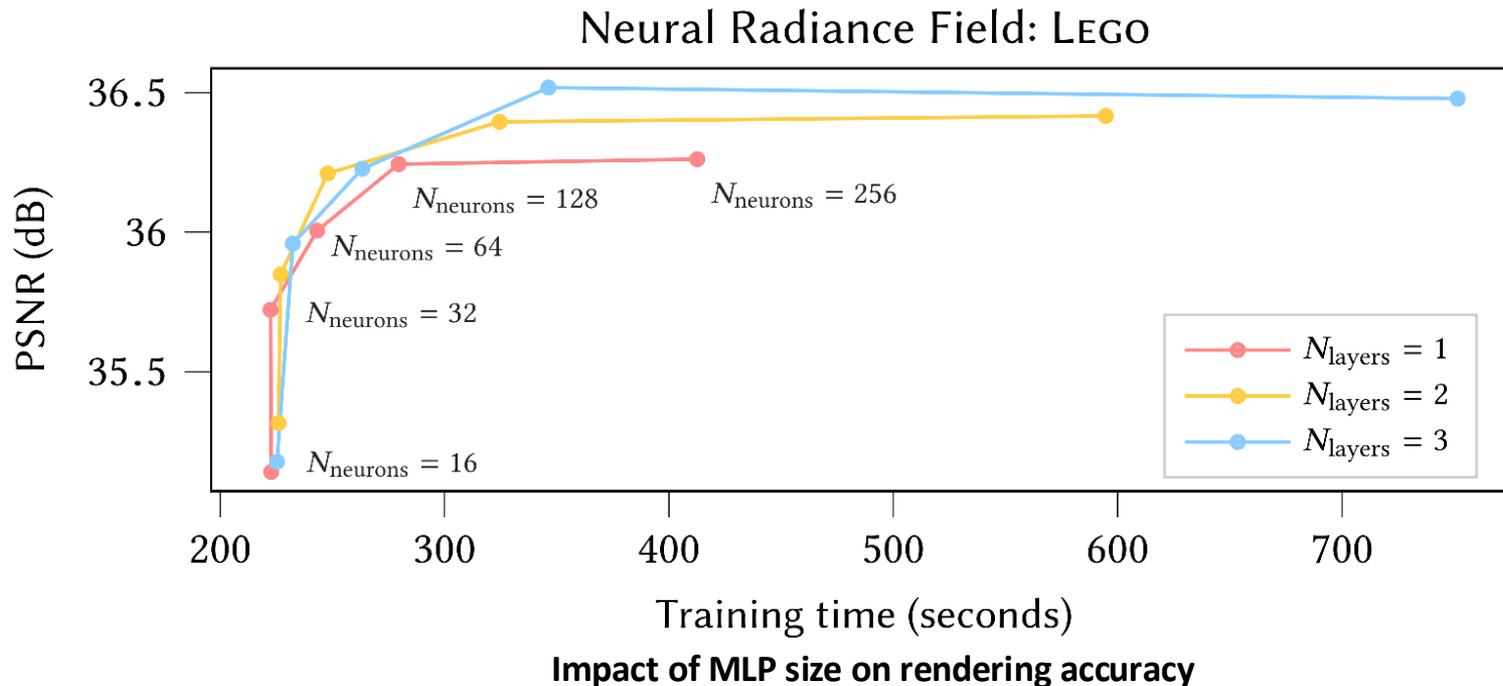
Outdoor free scene



Large indoor scene

## ➤ Challenge-1: Limitation of NeRF Model Capacity

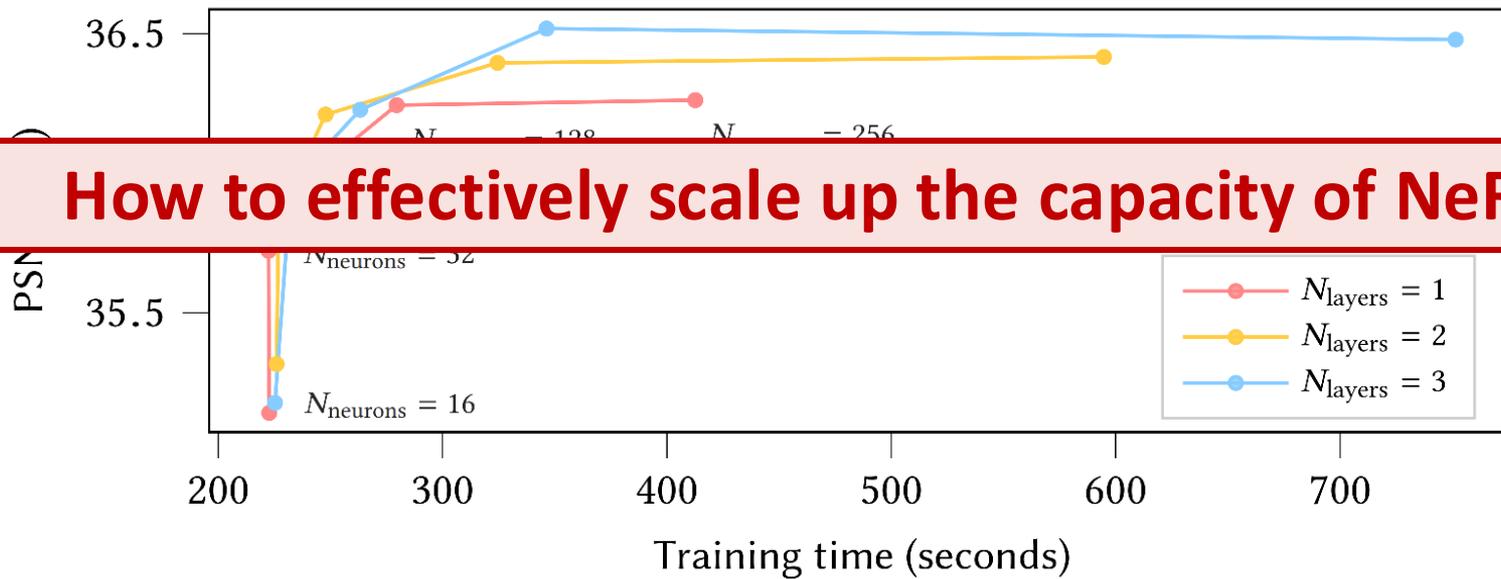
Directly increasing the network's size (width/depth) yields marginal performance improvement



## ➤ Challenge-1: Limitation of NeRF Model Capacity

Directly increasing the network's size (width/depth) yields marginal performance improvement

Neural Radiance Field: LEGO

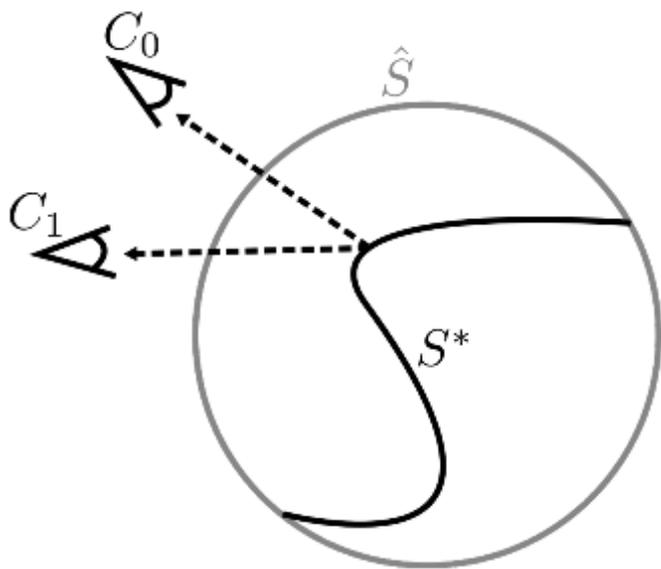


**How to effectively scale up the capacity of NeRF ?**

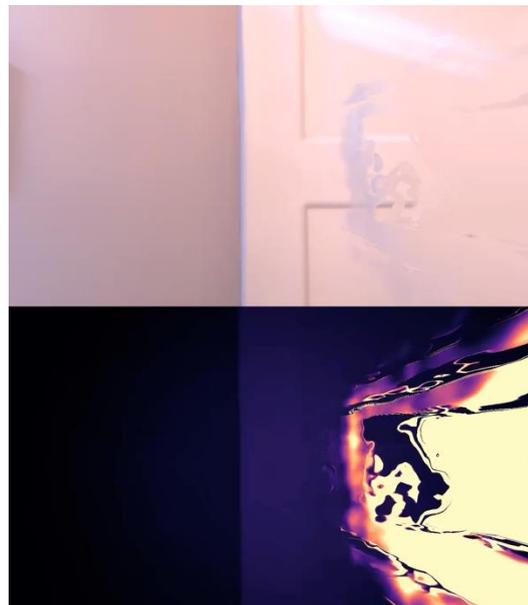
Impact of MLP size on rendering accuracy

## ➤ Challenge-2: Low Accuracy of Geometric Modeling

The quality of geometric modeling exhibits a significant influence on NeRF's generalization



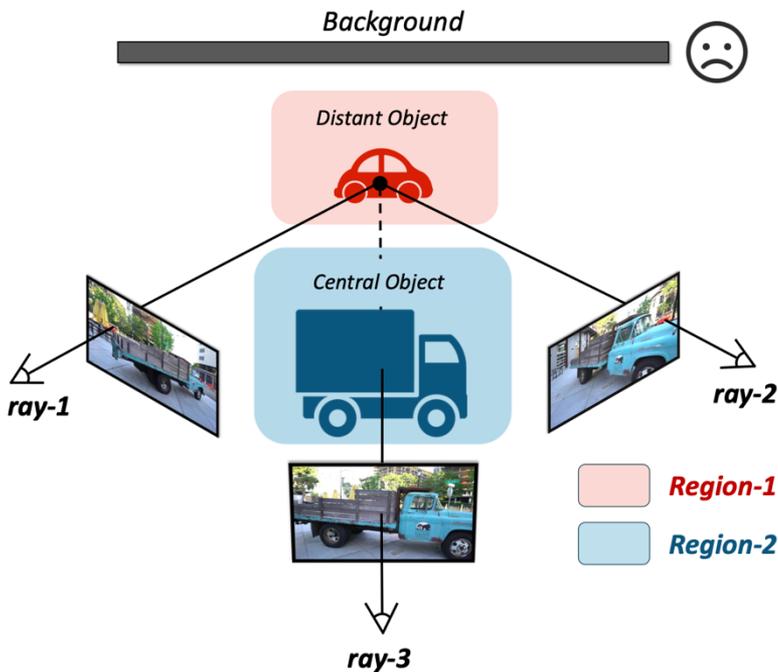
Shape-Radiance Ambiguity



# NeRF's Training interference



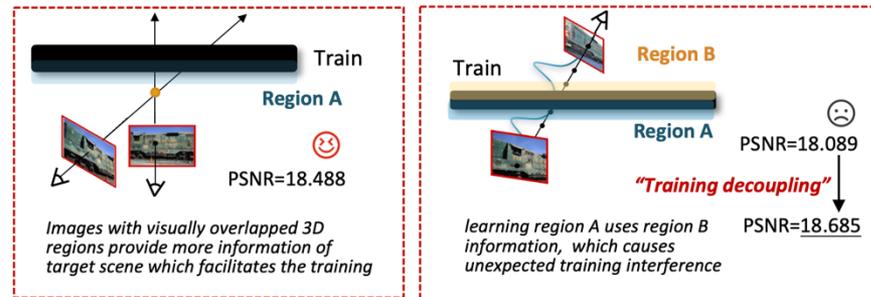
**Ray-3 does not contain valid information about distant object.**



The NeRF trained with more invisible rays (two sides of train) performs worse.



◀ Camera Pose    Front side A    Back side B



Without training interference

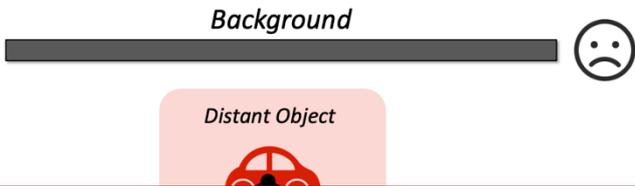
Decoupling mitigates training interference

# NeRF's Training interference

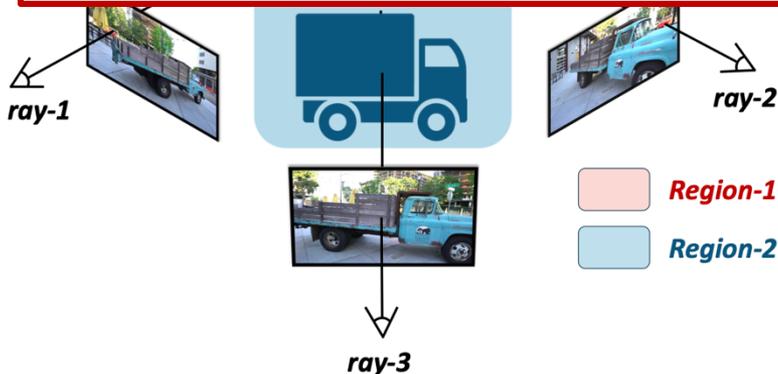


Ray-3 does not contain valid information about distant object.

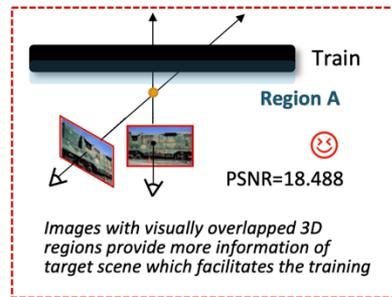
The NeRF trained with more invisible rays (two sides of train) performs worse.



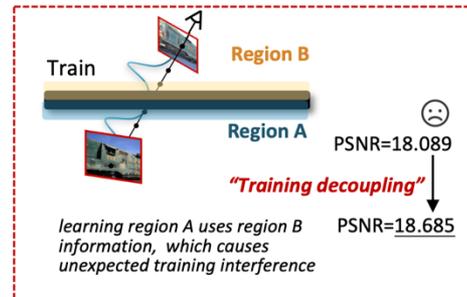
**Decouple the NeRF's Training in the Ray-dimension!**



◀ Camera Pose    ◻ Front side A    ◻ Back side B



Without training interference



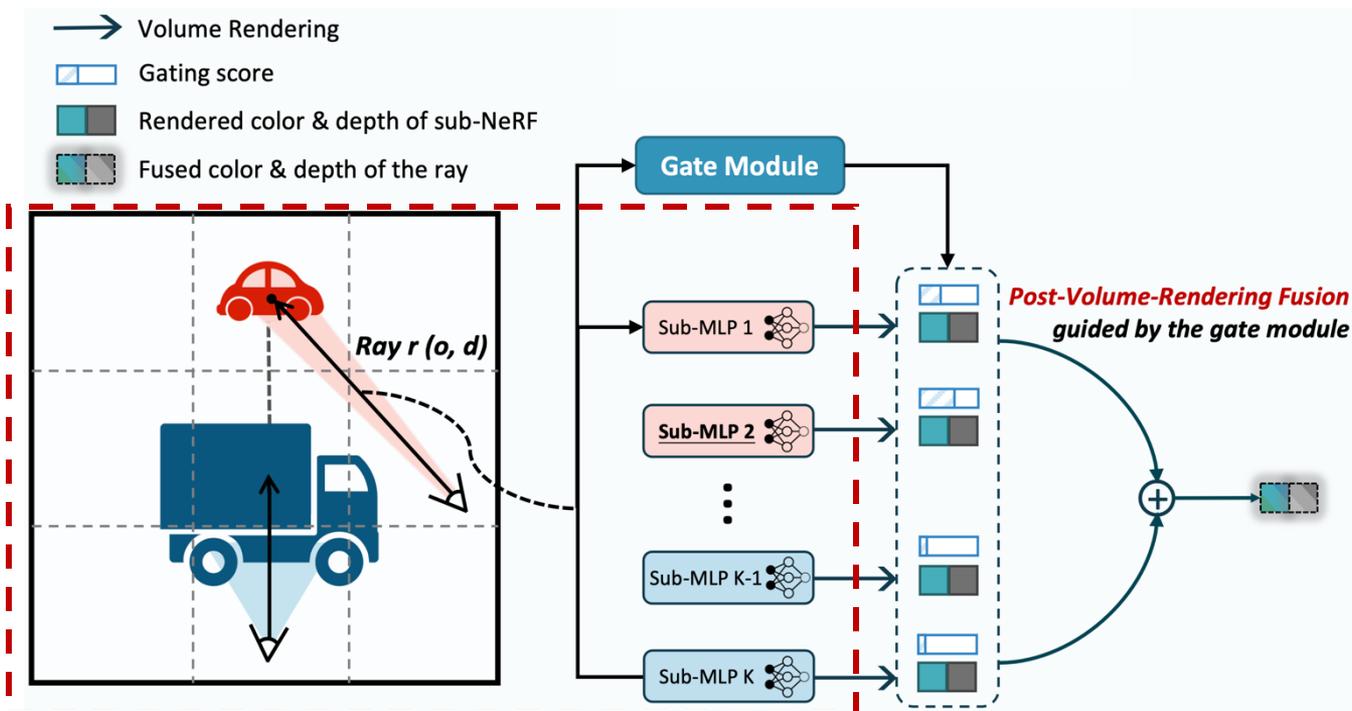
Decoupling mitigates training interference

Images with visually overlapped 3D regions provide more information of target scene which facilitates the training

learning region A uses region B information, which causes unexpected training interference

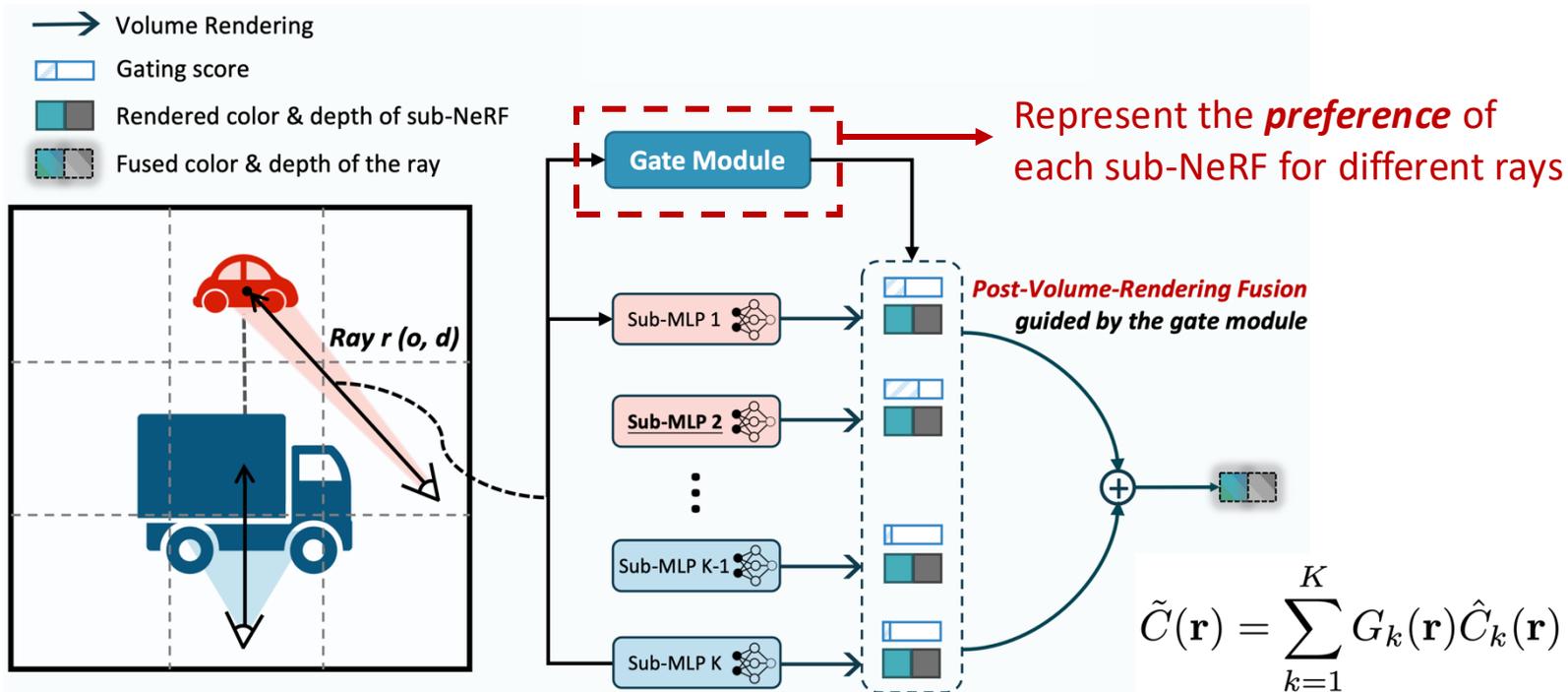
## Multi-NeRF Structure based on Hybrid Representation

- The feature grid is shared for all sub-NeRFs and the MLP decoders are independent



## Ray-wise Soft Gating Module Design

- A soft gating module is adopted to assign gating scores to the sub-NeRFs for each ray.



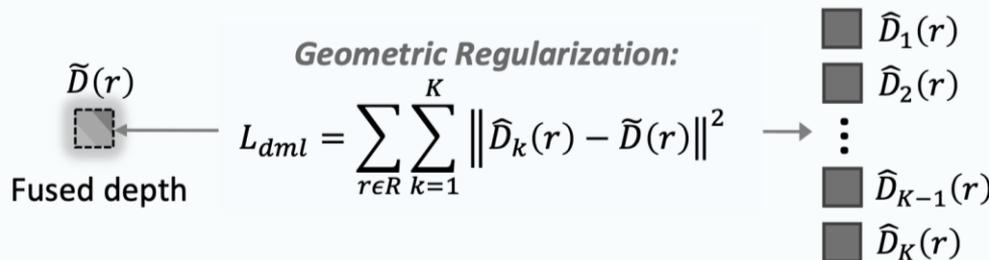
## ➤ Mutual Learning:

- each sub-NeRF not only learns from ground truth but also learns from each other.

**Sub-NeRFs *teach each other* with rendered depth to improve rendering consistency**

**Geometric Regularization:**

$$L_{dml} = \sum_{r \in R} \sum_{k=1}^K \|\hat{D}_k(r) - \tilde{D}(r)\|^2$$



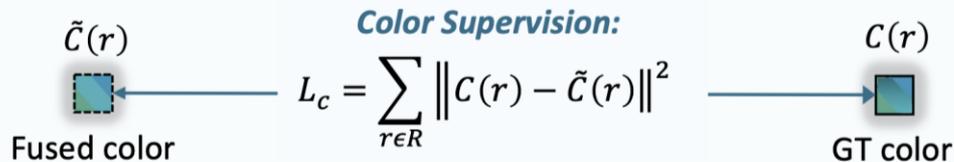
serves as geometric regularization



helps the model find more robust geometric solutions

**Color Supervision:**

$$L_c = \sum_{r \in R} \|C(r) - \tilde{C}(r)\|^2$$



Color rendering loss    Depth regularization    CV balancing loss

$$L = L_c + \lambda_1 L_{dml} + \lambda_2 L_{cv}$$

$$L_{cv} = \frac{\text{Var}(\overline{G}(\mathcal{R}))}{\left(\sum_{k=1}^n \overline{G}_k(\mathcal{R})/n\right)^2}$$

Encourages a balanced allocation of model parameters for training rays.

*prevents the gate module from collapsing onto a specific sub-NeRF*

$$\overline{G}_k(\mathcal{R}) = \sum_{\mathbf{r} \in \mathcal{R}} G_k(\mathbf{r}),$$

## ➤ Experiment Setup

- **Datasets: five datasets from different types of scenes**
  - (1) Object dataset: Masked Tanks-And-Temples (MaskTAT)
  - (2) 360-degree inward/outward-facing dataset:  
Tanks-And-Temples (TAT) & NeRF-360-v2 dataset
  - (3) Free-shooting-trajectory dataset:  
Free-dataset & ScanNet dataset
- **Baselines: different NeRF training frameworks**
  - (1) Grid-based single-NeRF: PlenOctrees, DVGO, Instant-NGP & F2-NeRF
  - (2) MLP-based single-NeRF: NeRF, NeRF++, MipNeRF & MipNeRF360
  - (3) Multi-NeRF frameworks: NGP-version of Block-NeRF, Switch-NeRF & Rad-NeRF

➤ Rad-NeRF achieves higher rendering quality than other single/multi-NeRF methods

Table 1: Quantitative results in complex scenes.

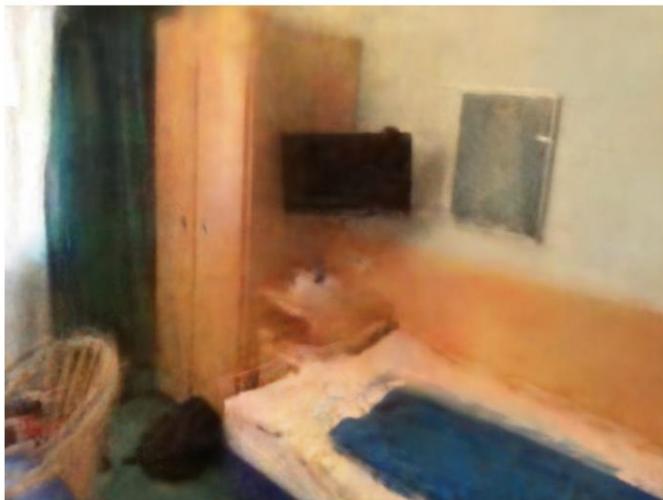
Methods	TAT			NeRF-360-v2			Free-Dataset		
	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓
NeRF++	20.419	0.663	0.451	27.211	0.728	0.344	24.592	0.648	0.467
MipNeRF360	<b>22.061</b>	<b>0.731</b>	<b>0.357</b>	<b>28.727</b>	<b>0.799</b>	<b>0.255</b>	<b>27.008</b>	<b>0.766</b>	0.295
MipNeRF360 <sub>short</sub> *	20.078	0.617	0.508	25.484	0.631	0.452	24.711	0.648	0.466
DVGO	19.750	0.634	0.498	25.543	0.679	0.380	23.485	0.633	0.479
Instant-NGP	20.722	0.657	0.417	27.309	0.756	0.316	25.951	0.711	0.312
F2-NeRF	–	–	–	26.393	0.746	0.361	26.320	<b>0.779</b>	<b>0.276</b>
Switch-NGP <sup>†</sup>	20.512	0.654	0.432	26.524	0.740	0.331	25.755	0.694	0.341
Block-NGP <sup>†</sup>	20.783	0.659	0.415	27.436	0.761	0.298	26.015	0.702	0.325
Rad-NeRF	<b>21.708</b>	<b>0.672</b>	<b>0.398</b>	<b>27.871</b>	<b>0.769</b>	<b>0.298</b>	<b>26.449</b>	0.719	<b>0.285</b>

\* MipNeRF360 requires nearly one day for training. For a fair comparison, we also report its results with one-hour of training.

† We adapt Switch-NeRF and Block-NeRF to the Instant-NGP fast training framework.

- Rad-NeRF achieves better recovery of scene details

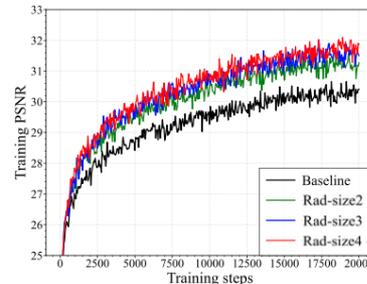
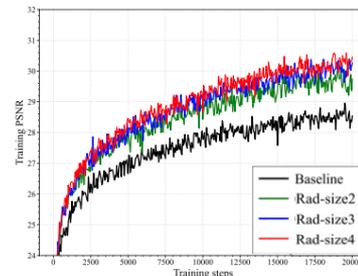
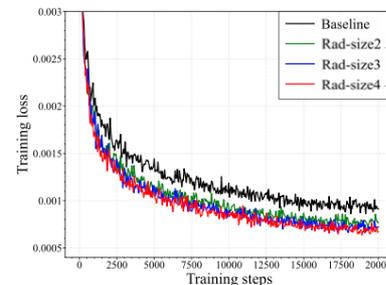
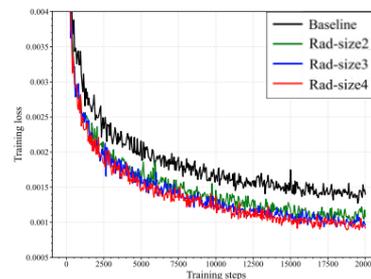
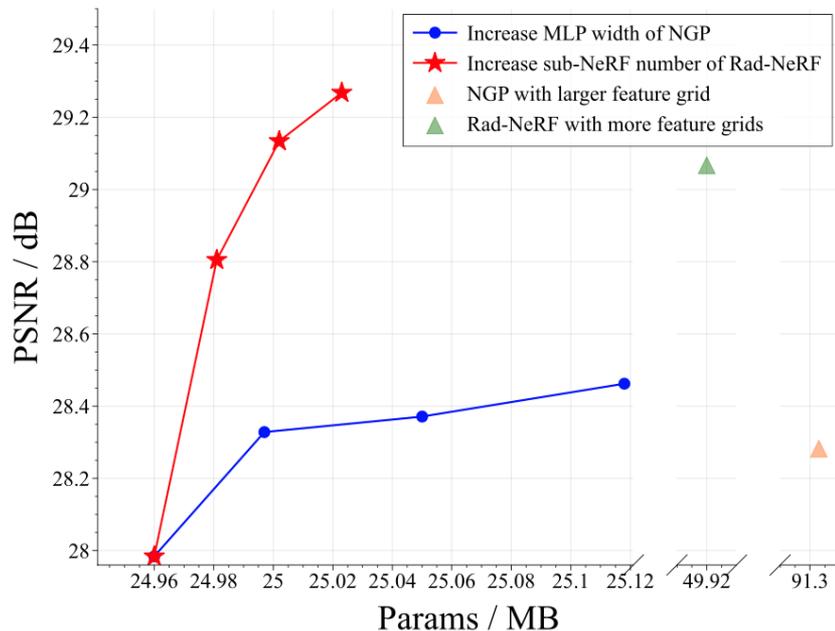
## Instant-NGP



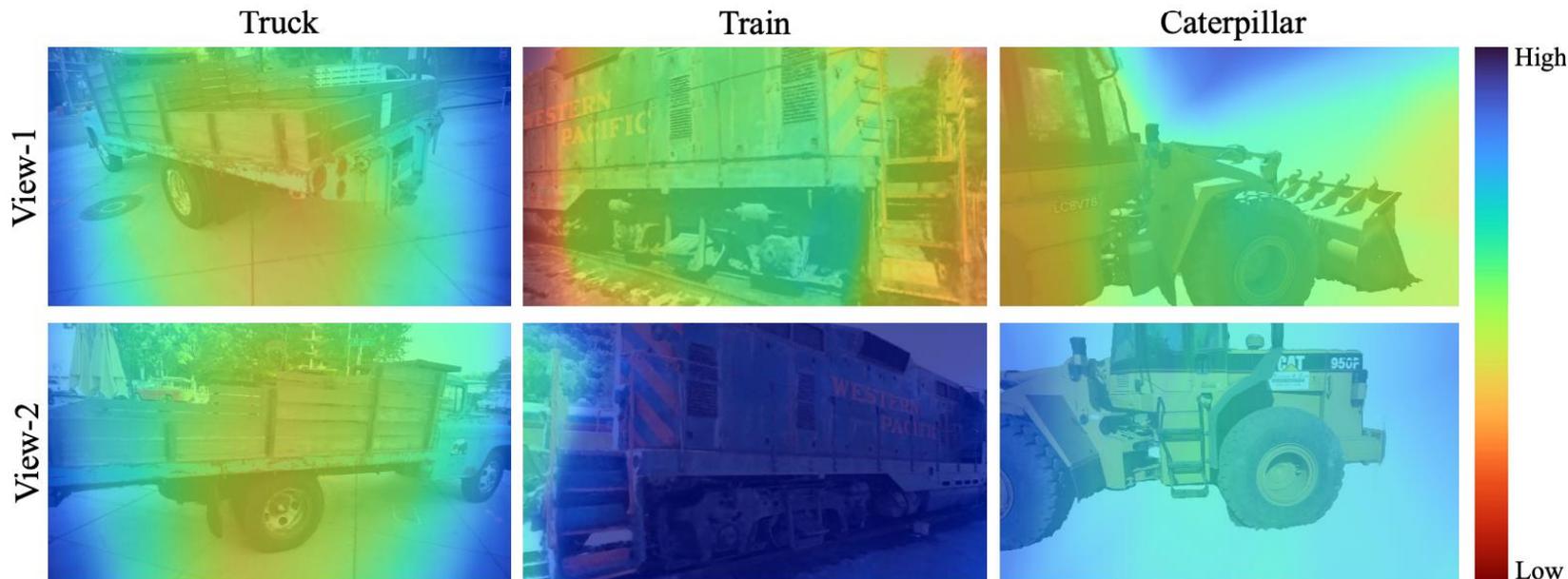
## Rad-NeRF



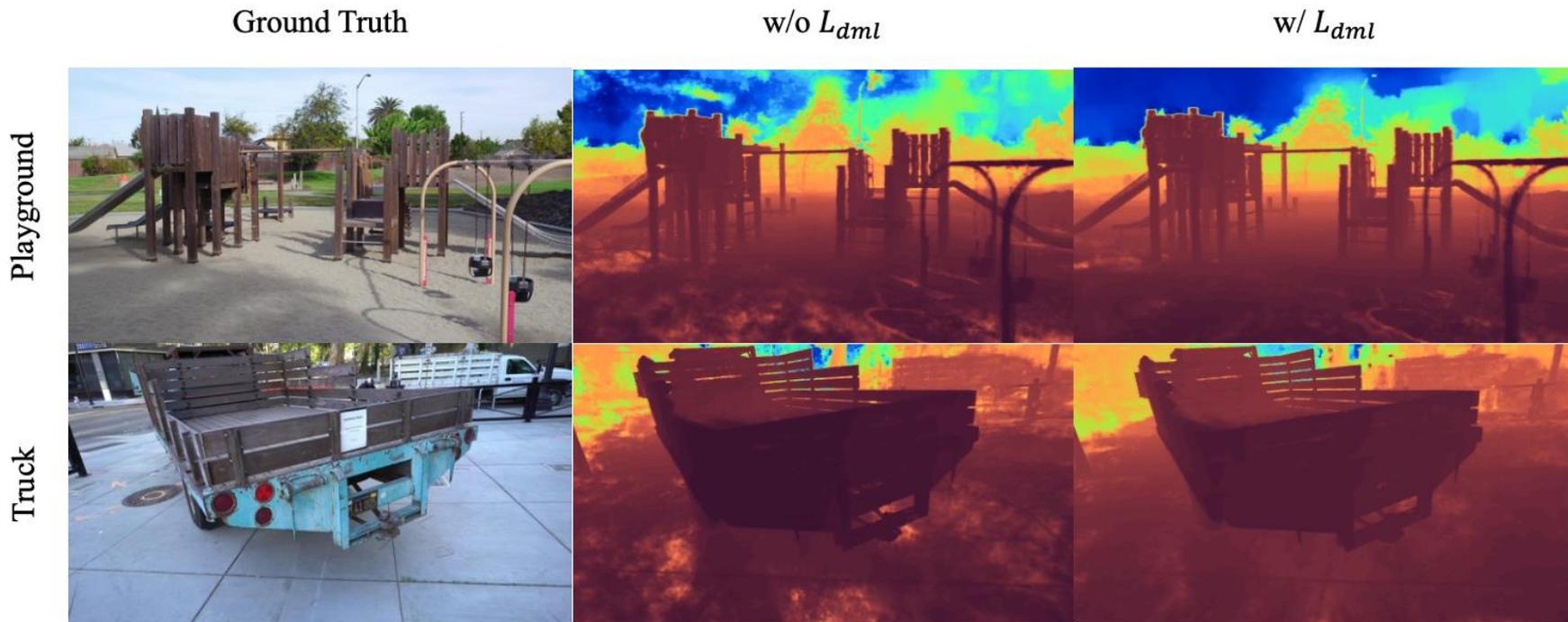
➤ Rad-NeRF achieves better **performance-parameter scalability**



- Rad-NeRF learns reasonable ray allocations, matching training interference “intuition”



- DML enables a smooth and reasonable depth prediction





清华大学电子工程系

Department of Electronic Engineering, Tsinghua University



美团



# Thanks for your attention!

## Q&A

Lidong Guo<sup>1\*</sup>, Xuefei Ning<sup>1\*</sup>, Yonggan Fu<sup>2</sup>, Tianchen Zhao<sup>3</sup>, Zhuoliang Kang<sup>3</sup>,  
Jincheng Yu<sup>1</sup>, Yingyan (Celine) Lin<sup>2</sup>, Yu Wang<sup>1</sup>

<sup>1</sup>Tsinghua University, <sup>2</sup>Georgia Institute of Technology, <sup>3</sup>Meituan

E-mail: [gld21@mails.tsinghua.edu.cn](mailto:gld21@mails.tsinghua.edu.cn), [foxdoraame@gmail.com](mailto:foxdoraame@gmail.com), [yu-wang@tsinghua.edu.cn](mailto:yu-wang@tsinghua.edu.cn)

