

# Recent Advances in Efficient Neural Light Field

Huan Wang

Northeastern University, Boston, USA

Talk @Zhidx, 07/07/2023, Fri

**Short Bio:** Huan Wang, BE'16@ZJU, MS'19 @ZJU, now 4th-year Ph.D. candidate @NEU. Interned at Alibaba/MERL/Snap. Now interning at Google. Work on **efficient deep learning** (network sparsity & distillation, etc.), primarily in the vision domain.

## Outline: 2 Papers

- R2L: Distilling NeRF to NeLF (ECCV 2022)
- MobileR2L: Run R2L on mobile devices via arch. optimization (CVPR 2023)

# R2L: Distilling Neural Radiance Field to Neural Light Field for Efficient Novel View Synthesis (ECCV 2022)

Huan Wang<sup>1,2,\*</sup> Jian Ren<sup>1,†</sup> Zeng Huang<sup>1,‡</sup> Kyle Olszewski<sup>1</sup> Menglei Chai<sup>1</sup> Yun Fu<sup>2</sup>  
and Sergey Tulyakov<sup>1</sup>

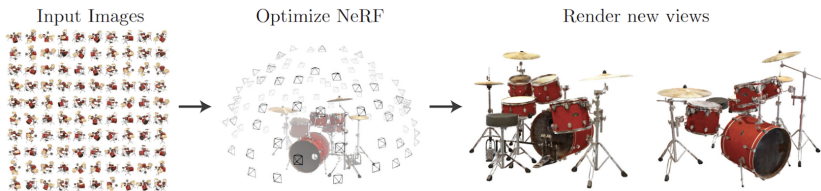
<sup>1</sup>Snap Inc. <sup>2</sup>Northeastern University  
Project: <https://snap-research.github.io/R2L>

Talk @Zhidx, 07/07/2023, Fri



# Background & Motivation

# What is NVS (Image-Based Rendering) and NeRF?



Take a few photos  
around an object



Train a neural  
network model



How does the object  
look like from a *new*  
perspective?

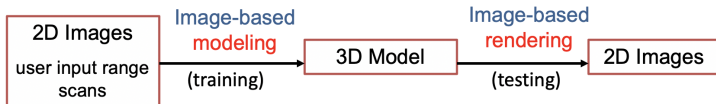


Figure: Illustration of NeRF framework. Src: ([Mildenhall et al., 2020](#)); edited.

# NeRF, Implicit Scene Representation

- NeRF (or neural radiance field) ([Mildenhall et al., 2020](#)) is an *implicit* representation of a scene – a buzz word since 2020, Best Paper HM in ECCV'20.
- It opens the new doors of representing complex scenes with a simple MLP network.
- So simple but working!  $(x, y, z; \theta, \phi) \rightarrow (\text{NeRF network}) \rightarrow (r, g, b; \sigma) \rightarrow (\text{Accumulation}) \rightarrow (r, g, b)$

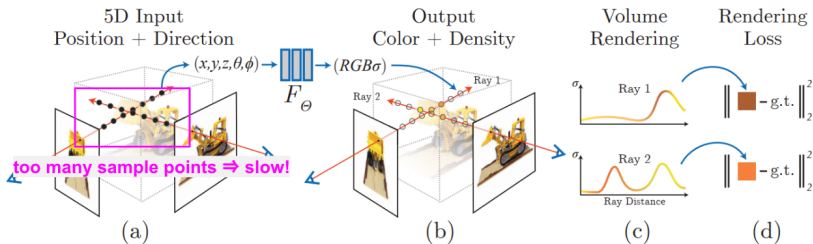


Figure: Illustration of NeRF framework. Src: ([Mildenhall et al., 2020](#)); edited.

# Aside, Background: 3D Representations

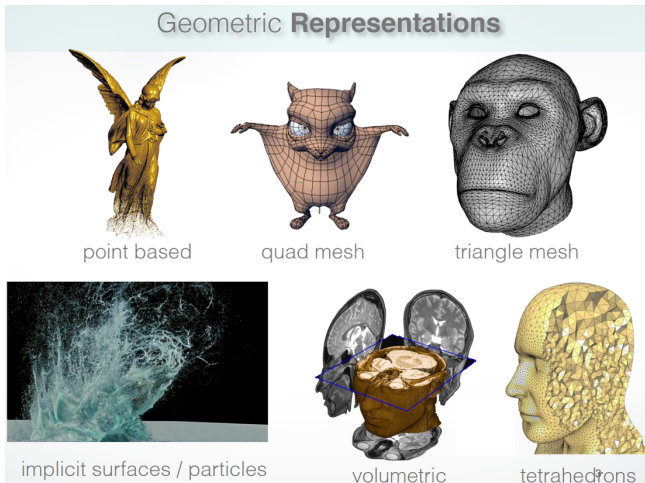


Figure: Courtesy: Prof. Hao Li, USC, CSCI 621: Digital Geometry Processing, <http://cs621.hao-li.com>. *Different representations require different processing methods.*

# Aside, Background: Implicit Representations

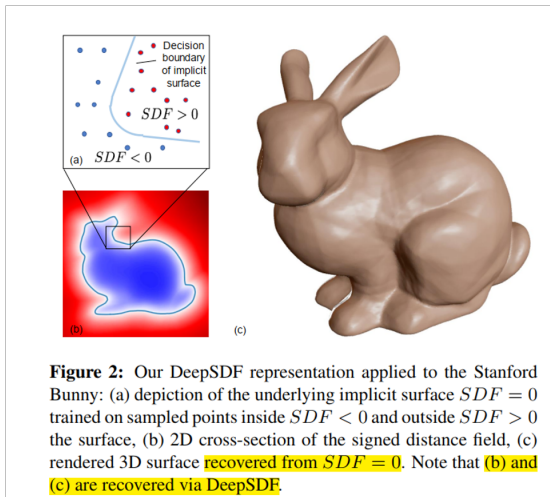


Figure: DeepSDF (Park et al., 2019) uses SDF (signed distance function) parameterized by a neural network to represent the surface / shape of an object.



# Motivation: NeRF is great but *slow* in inference

- However, one primary downside of NeRF is the **prohibitively slow** inference: Rendering one  $400 \times 400$  image takes 6.7s on an NVIDIA V100 GPU.
- **Our goal:** *Faster inference* of NeRF for efficient NVS (novel view synthesis)

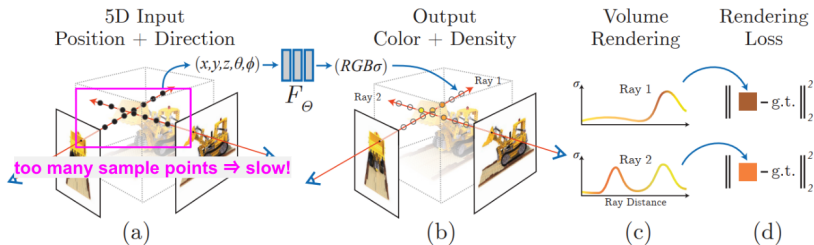


Figure: Illustration of NeRF framework. Src: ([Mildenhall et al., 2020](#)); edited.

# Speed Bottleneck: Too Many Sampled Points

- NeRF is slow because of too many sampled points along the ray. Tradeoff here: quality vs. speed.
- **Our goal: *Faster inference* of NeRF for efficient NVS (novel view synthesis)**  
 ⇒ Technical goal: **Reduce the #sampled points**

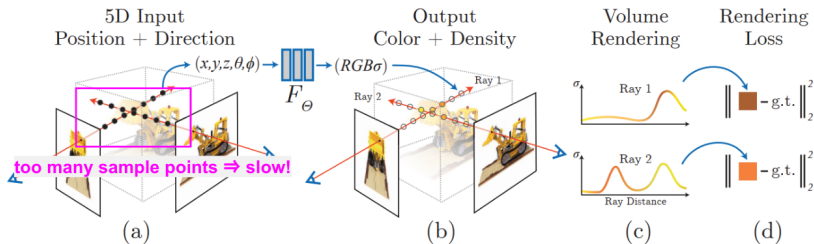
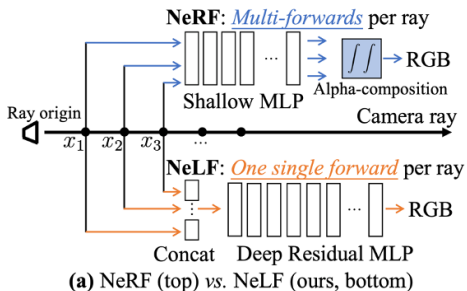


Figure: Illustration of NeRF framework. Src: ([Mildenhall et al., 2020](#)); edited.

## Proposed Method: R2L (NeRF to NeLF)

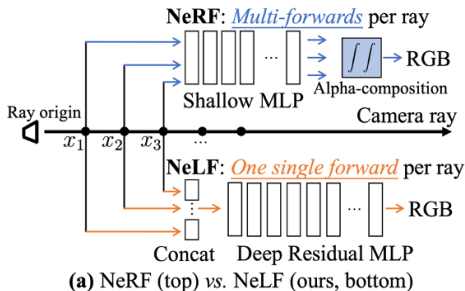
# NeLF vs. NeRF

- We seek fast inference through another representation of a scene: **NeLF (or neural light field)**.
- NeLF vs. NeRF – **Upside**: Fast! Rendering one pixel amounts to **one** network query in NeLF vs. **hundreds** of network queries in NeRF, as shown below.



# NeLF vs. NeRF

- We seek fast inference through another representation of a scene: **NeLF (or neural light field)**.
- NeLF vs. NeRF – **Upside**: Fast! Rendering one pixel amounts to **one** network query in NeLF vs. **hundreds** of network queries in NeRF, as shown below.



**Downsides:** (1) NeLF is inherently **harder** to learn than NeRF. (2) Given the same batch of images, the data for NeLF is **much fewer** than NeRF.

In short, to learn NeLF successfully, we must overcome **two problems**: #1 The target function is harder to learn. #2 Even worse, the data sample size shrinks.

# R2L: ResMLP

To overcome problem #1, we make two contributions: (1) ResMLP, (2) Better ray representation.

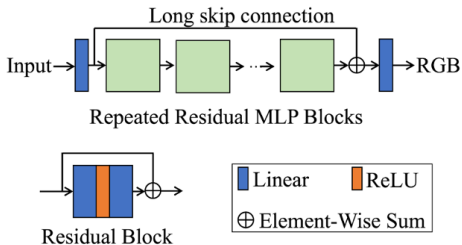


Figure: Architecture of our proposed ResMLP network. Note, this is the **first attempt** that deep residuals are used in NeRF-like models.

# R2L: Ray Representation - PointCat

To overcome problem #1, we make two contributions: (1) ResMLP, (2) Better ray representation.

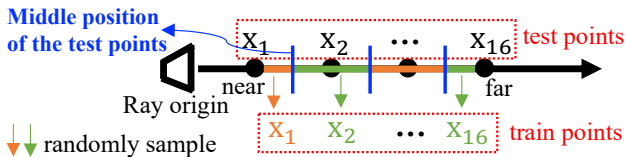
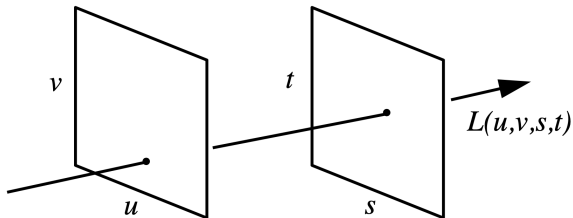


Figure: Illustration of the point sampling in training and testing of our method. The orange and green colors denote the different *segments* of the ray. The blue color marks the *start* and *end* points of each segment. Each sampled train point is colored *based on the corresponding segment color*.

The sampled points are **concatenated** into a long vector to make the input of our NeLF network – very simple!

During training, the points are **randomly sampled**; during testing, they are **fixed**.

# Other Ray Representations - Two Planes



**Figure 1:** The light slab representation.

Figure: Two-plane parameterization (Levoy & Hanrahan, 1996) represents each ray in a 3D space with its intersection points on two planes.



## Other Ray Representations - Plücker

mathematically convenient. In particular, we propose to leverage the **6D Plücker parameterization** of the space of light rays  $\mathcal{L}$  for LFNs. The Plücker coordinates (see [49] for an excellent overview) of a ray  $\ell$  through a point  $\mathbf{p}$  in a normalized direction  $\mathbf{d}$  are

$$\mathbf{r} = (\mathbf{d}, \mathbf{m}) \in \mathbb{R}^6 \text{ where } \mathbf{m} = \mathbf{p} \times \mathbf{d}, \text{ for } \mathbf{d} \in \mathbb{S}^2, \mathbf{p} \in \mathbb{R}^3. \quad (4)$$

where  $\times$  denotes the **cross product**. While Plücker coordinates are **a-priori 6-tuples of real numbers**, the coordinates of any ray lie on a 4-dimensional subspace  $\mathcal{L}$ . Plücker coordinates uniformly represent

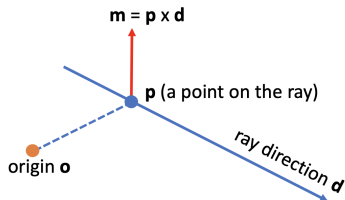


Figure: Plücker representation (Sitzmann et al., 2021) is a tuple of 6 real numbers, made with a normalized direction vector (3D) and a cross-product (3D) of the direction and a point on the line. **Not as good as our Point-Concat scheme** (see ablation study in the paper).

# R2L: Distilling w/ Pseudo Data & Finetuning w/ Real

**Network larger, more data hungry?** (i.e., problem #2)  $\Rightarrow$  Use a *pretrained* NeRF to synthesize a large amount of data.

**Training loss:** MSE

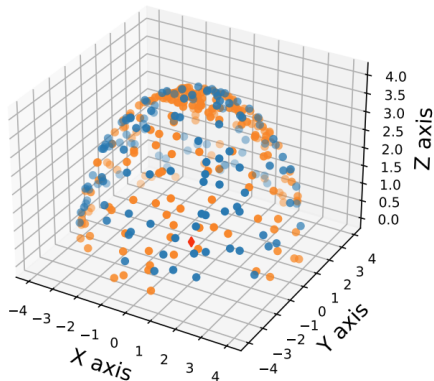
$$\mathcal{L} = \text{MSE} \left( G_{\phi} \left( \underbrace{(x_o, y_o, z_o, x_d, y_d, z_d)}_{\text{Pseudo data (via querying a trained NeRF)}} \right), (\hat{r}, \hat{g}, \hat{b}) \right)$$

R2L network    Ray origin    Ray direction    Ray RGB

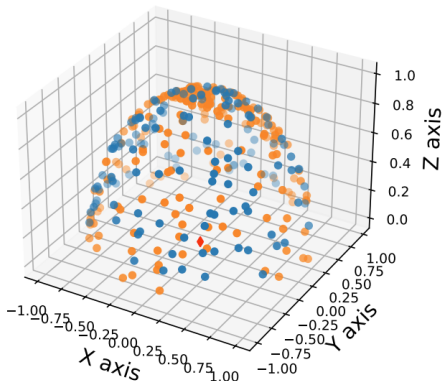
**Finetune:** After training our ResMLP model on the pseudo data, it can achieve comparable performance to the NeRF teacher. Then we *finetune* the student on the original real data, which *significantly* boosts the student's performance.

# Final Question: How to Get the Psuedo Inputs?

Sampling from joint distribution  $\Rightarrow$  marginal distribution product (*i.e.*, assuming independent):  $\mathcal{D}(X, Y, Z) \approx \mathcal{D}(X) \times \mathcal{D}(Y) \times \mathcal{D}(Z)$



(a) Origins

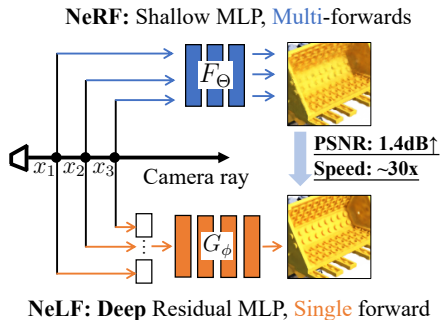


(b) Directions

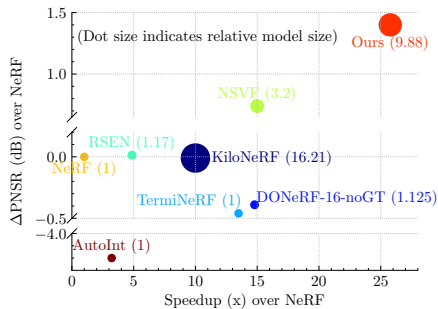
Figure: **Blue**: original data. **Orange**: pseudo data. Pseudo data needs to overlap with the original data – same domain!

# Experimental Results

# Comparison Overview (NeRF Synthetic Dataset)



(a) NeRF vs. our NeLF method



(b) Speedup-PSNR-Model Size comparison

Figure: (a) Our neural light field (NeLF, bottom) method **improves** the rendering quality by **1.40 PSNR** over neural radiance field (NeRF, top) (Mildenhall et al., 2020) on the NeRF synthetic dataset, while being around **30× faster**. (b) Our method achieves a more favorable speedup-PSNR-model size tradeoff than other efficient novel view synthesis methods on the NeRF synthetic dataset.

# Visual Comparison (NeRF Synthetic Dataset)

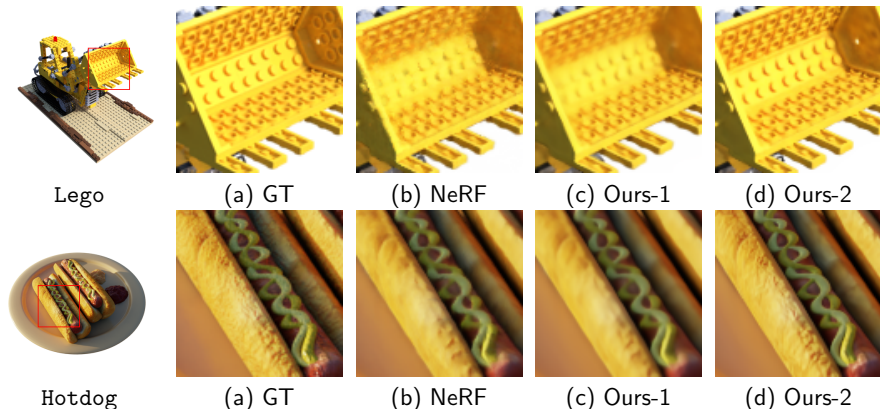


Figure: Ours vs. NeRF (Mildenhall et al., 2020). Ours-1 is trained solely on pseudo data, Ours-2 on pseudo and real data. **Much better quality** with only **1/26 FLOPs**. Please see our paper for more results.

# Visual Comparison (Real-World Dataset)

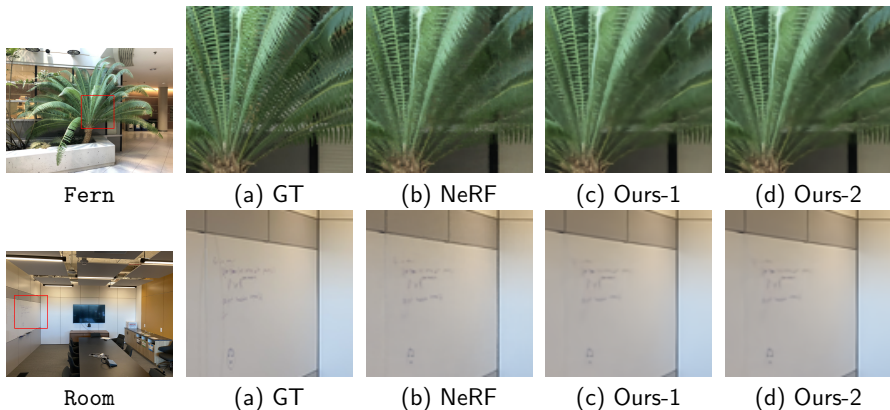


Figure: Ours vs. NeRF (Mildenhall et al., 2020). Ours-1 is trained solely on pseudo data, Ours-2 on pseudo and real data. **Comparable quality** with only **1/26 FLOPs**. Please see our paper for more results.

# Conclusion



# Conclusion

We present **R2L**, a new and novel **efficient NVS** method by **distilling a NeRF to NeLF**. Three major techniques are proposed to enable it:

- **Network**: We propose a **much deeper ResMLP** network to seek greater representative power – This is the *first attempt* that optimizes the NeRF speed from the network architecture level and the *first* time dense residuals are properly integrated into a NeRF network.
- **Data**: We propose to employ pretrained NeRF to **synthesize abundant pseudo data**. The model is later **finetuned** on the real data to further boost performance.
- **Ray representation**: We represent the ray with a **simple concat scheme** and propose to **add noise** during training to curb overfitting.

Empirically, to our knowledge, this is the *first NeLF* that runs round **30× faster** than NeRF while still being **significantly better**.

# GitHub Code & Limitations

Thank you!

Code: <https://github.com/snap-research/R2L>

To Improve:

- **Slow training**
- **Still not real-time on mobile devices**  $\Rightarrow$  MobileR2L
- Needs a pretrained teacher? (Can be an upside too! Better teacher, better student)
- Better random sampling (accelerate training)
- ...

# Real-Time Neural Light Field on Mobile Devices (CVPR 2023)

Junli Cao<sup>1</sup> Huan Wang<sup>2</sup> Pavlo Chemerys<sup>1</sup> Vladislav Shakhrai<sup>1</sup> Ju Hu<sup>1</sup>  
Yun Fu<sup>2</sup> Denys Makoviichuk<sup>1</sup> Sergey Tulyakov<sup>1</sup> Jian Ren<sup>1</sup>

<sup>1</sup>Snap Inc. <sup>2</sup>Northeastern University

Project: <https://snap-research.github.io/MobileR2L>

Talk @Zhidx, 07/07/2023, Fri



# Motivation

Run R2L in real-time on mobile devices? Problems:

- Slow! 3 seconds for rendering 200x200 on iPhone 13.
- $800 \times 800$  input,  $800 \times 800$  output? Too large, OOM!

# Learn a Page from Image Super-Resolution (SR)

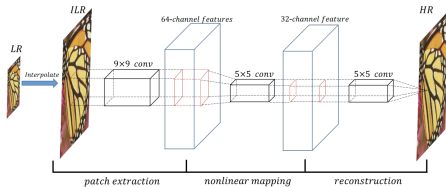
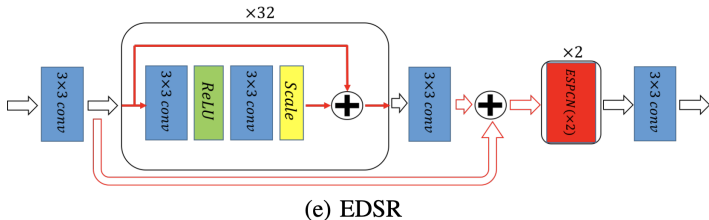


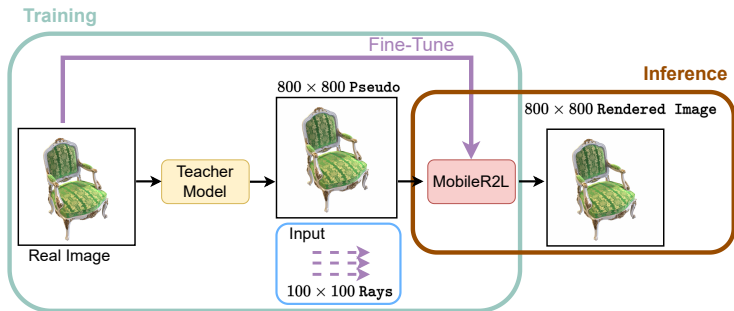
Figure 2: Sketch of the SRCNN architecture.



(e) EDSR

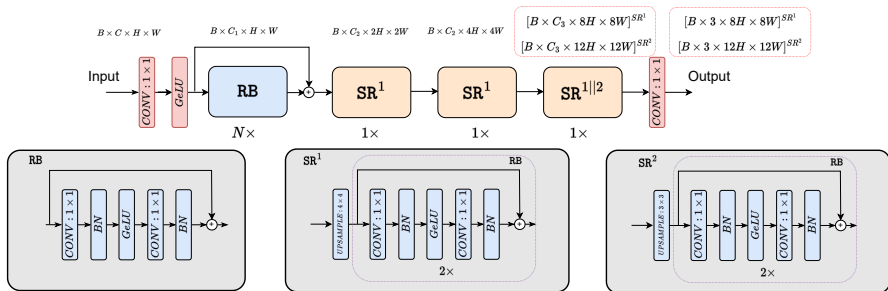
Figure: Illustration of SRCNN (Dong et al., 2014) (ECCV'14) and EDSR (Lim et al., 2017) (CVPRw'17). Src: (Yang et al., 2019). Want to be faster? Upsampling later.

# Overview of MobileR2L



**Figure: Training and Inference Pipeline.** The training involves a teacher model to generate pseudo data, which is used to learn the MobileR2L. The teacher model, e.g., NeRF, is trained on real images. Once we have the teacher model, we use it to generate pseudo images, e.g., images with the resolution of  $800 \times 800$ , in addition to down-scaled rays, e.g., rays with spatial size as  $100 \times 100$ , that share the same origin with the pseudo images to train the MobileR2L. After that, we use the real data to fine-tune MobileR2L. For inference, we directly forward the rays into the pre-trained MobileR2L to render images.

# The Proposed Architecture of MobileR2L



**Figure: Overview of Network.** The backbone includes residual blocks (RB) repeated that are 28 times ( $N = 28$ ). Following the backbone, there are two types of super-resolution (SR) modules. The first SR module ( $\text{SR}^1$ ) has kernel size  $4 \times 4$  in the Transpose CONV layer that doubles the input  $H, W$  to  $2H, 2W$ , whereas the second SR module ( $\text{SR}^2$ ) has kernel size  $3 \times 3$ , tripling the spatial size to  $3H, 3W$ .

# Other Changes from R2L to MobileR2L

- Input Format:
  - R2L: Use MLP. Rays are mutually independent. No norm layers  $\Rightarrow$  **Shuffled rays as input**
  - MobileR2L: Use Conv. Rays are not mutually independent  $\Rightarrow$  **Cropped patch as input** (use BN for training stability)
- 1x1 MLP to Conv (27% faster! see below)
- ReLU to GeLU (PSNR slightly better, a minor change)

Table 4. **Analysis of Network Design.** For all the comparisons, we use the input tensor with the spatial size as  $100 \times 100$  and render the image with spatial size. The latency (ms) is measured on iPhone 13 (iOS16) with models compiled with CoreMLTools [11]

	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Latency $\downarrow$
MLP	19.13	0.9759	0.6630	19.57
CONV2D	19.16	0.9759	0.6301	14.30
CONV2D + ReLU	26.82	0.9949	0.0282	16.20
CONV2D + GeLU	26.99	0.9949	0.0730	17.00
CONV2D + GeLU + BN	<b>27.18</b>	<b>0.9954</b>	<b>0.0259</b>	17.00



# Experimental Results

# Quantitative Comparison

Table: **Quantitative Comparison** on Synthetic 360° and Forward-facing. Our method obtains better results on the three metrics than NeRF for the two datasets. Compared with MoibleNeRF and SNeRG, we achieve better results on most of the metrics.

	Synthetic 360°			Forward-facing		
	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓
NeRF ( <a href="#">Mildenhall et al., 2020</a> )	31.01	0.947	0.081	26.50	0.811	0.250
NeRF-Pytorch ( <a href="#">Yen-Chen, 2020</a> )	30.92	0.991	0.045	26.26	0.965	0.153
SNeRG ( <a href="#">Hedman et al., 2021</a> )	30.38	0.950	<b>0.050</b>	25.63	0.818	<b>0.183</b>
MoibleNeRF ( <a href="#">Chen et al., 2023</a> )	30.90	0.947	0.062	25.91	0.825	<b>0.183</b>
MobileR2L (Ours)	<b>31.34</b>	<b>0.993</b>	0.051	<b>26.15</b>	<b>0.966</b>	0.187
Our Teacher	33.09	0.961	0.052	26.85	0.827	0.226

Our Teacher here: Mip-NeRF ([Barron et al., 2021](#)) (Better teacher, better student).

# Visual Comparison

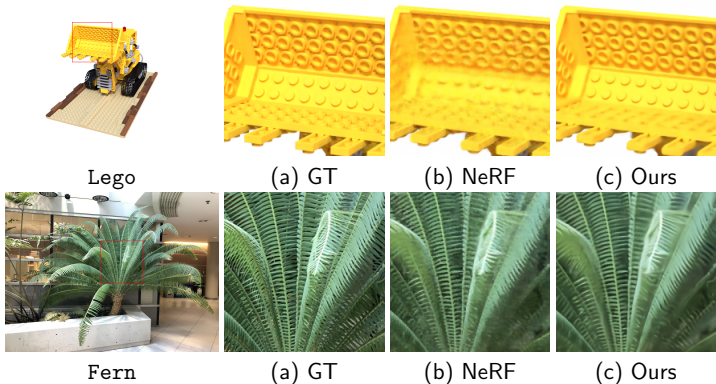


Figure: Visual comparison between our method and NeRF (Mildenhall et al., 2020) (trained via NeRF-Pytorch (Yen-Chen, 2020)) on the synthetic 360° Lego (size:  $800 \times 800 \times 3$ ) and real-world forward-facing scene Fern (size:  $1008 \times 756 \times 3$ ). *Best viewed in color.* Please refer to our [webpage](#) for more visual comparison results.

# Zoom-in Comparison



Figure: Zoom-in comparisons. *Top row*: MobileNeRF (Chen et al., 2023). Results are obtained from the code and demo released by the authors. *Bottom row*: MobileR2L. Our approach renders high-quality images even for zoom-in views.

# Storage (Representation Size) Comparison

Table: **Analysis of Storage (MB)** required for different rendering methods. Our method has a clear advantage over existing works with much less storage required, compared to MobileNeRF ([Chen et al., 2023](#)) and SNeRG ([Hedman et al., 2021](#)).

	Synthetic 360°			Forward-facing		
	MobileNeRF	SNeRG	Ours	MobileNeRF	SNeRG	Ours
Disk storage	125.8	86.8	<b>8.3</b>	201.5	337.3	<b>8.3</b>

**15× ~ 24× smaller!**

# Latency Comparison on Mobile Devices

**Table: Analysis of Inference Speed.** Latency (ms) is obtained on iPhone with iOS 16. Following MobileNeRF (Chen et al., 2023), we use the notation  $\frac{M}{N}$  to indicate that  $M$  out of  $N$  scenes in the Forward-facing dataset that can not run on devices. Specifically, **MobileNeRF cannot render Leaves and Orchids in Forward-facing.**

	Synthetic 360°		Forward-facing	
	MobileNeRF	Ours	MobileNeRF	Ours
iPhone 13	17.54	26.21	27.15 $\frac{2}{8}$	18.04
iPhone 14	16.67	22.65	20.98 $\frac{2}{8}$	16.48

Note, MobileNeRF is 1/12 Award Candidates in CVPR'23. MobileR2L presents an equally (if not more) promising alternative for efficient mobile neural rendering.

## More Demo: Snap's AR Demo

- A fake lego bulldozer on a real desk.

## Another Remark: Why MobileR2L is Special?

- It is the *first* to propose a **2D-3D hybrid architecture** design for NVS, opening the new doors to marry the two areas together – just image so many techniques from the 2D image restoration area; many of them can be used here for improved quality in NeRF/NeLF.



## Conclusion and Discussions

# Conclusion

- We present MobileR2L, the very first NeLF that can run in real-time on mobile devices.
- MobileR2L is featured by a hybrid design of 3D light field network and 2D super-resolution architecture, enabling low-resolution ( $100 \times 100$ ) input with high-resolution ( $800 \times 800$ ) output.
- Empirically, MobileR2L achieves  $\sim$ **60fps** real-time rendering on iPhone14, while saving over  $24\times$  storage than MobileNeRF (one of the **12 Award Candidate papers** in CVPR 2023).

# Limitations / Future Work

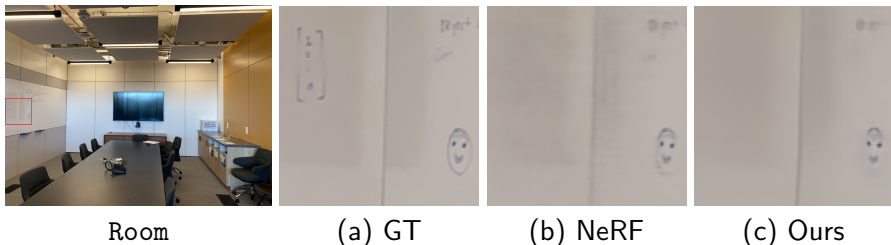


Figure: Visual comparison on the real-world scene Room. Both our model and NeRF fail to synthesize the whiteboard writings on the upper-left of the cutout patch.

- **Quality.** On some scenes, we still observe blurry rendered results, fine-grained details and textures are missing (see above).
- **Training speed.** R2L or MobileR2L training cost is  $100 \sim 500\times$  of NeRF (so essentially, R2L or MobileR2L is trading inference cost with training cost).
- **Still need a teacher.** How to learn light field directly remains a question.

## Thanks! Questions?

More materials:

- CVPR'23 Workshop: *Efficient Neural Networks: From Algorithm Design to Practical Mobile Deployments*, hosted by Snap. [[Webpage](#)]
- Slides can be found on my personal [webpage](#).

Feel free to drop me a line ([wang.huan@northeastern.edu](mailto:wang.huan@northeastern.edu)) or raise a GitHub issue if you have any questions.

# References

- Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *ICCV*, 2021.
- Zhiqin Chen, Thomas Funkhouser, Peter Hedman, and Andrea Tagliasacchi. Mobilenerf: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures. In *CVPR*, 2023.
- Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional network for image super-resolution. In *ECCV*, 2014.
- Peter Hedman, Pratul P Srinivasan, Ben Mildenhall, Jonathan T Barron, and Paul Debevec. Baking neural radiance fields for real-time view synthesis. In *ICCV*, 2021.
- Marc Levoy and Pat Hanrahan. Light field rendering. In *Proceedings of the Annual Conference on Computer Graphics and Interactive Techniques*, 1996.
- Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *CVPR Workshop*, 2017.
- Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.
- Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *CVPR*, 2019.