

Learning Neural Duplex Radiance Fields for Real-Time View Synthesis

Ziyu Wan¹ Christian Richardt² Aljaž Božič² Chao Li² Vijay Rengarajan²
Seonghyeon Nam² Xiaoyu Xiang² Tuotuo Li² Bo Zhu² Rakesh Ranjan² Jing Liao^{1*}
¹City University of Hong Kong ²Meta Reality Labs
raywzy.com/NDRF/

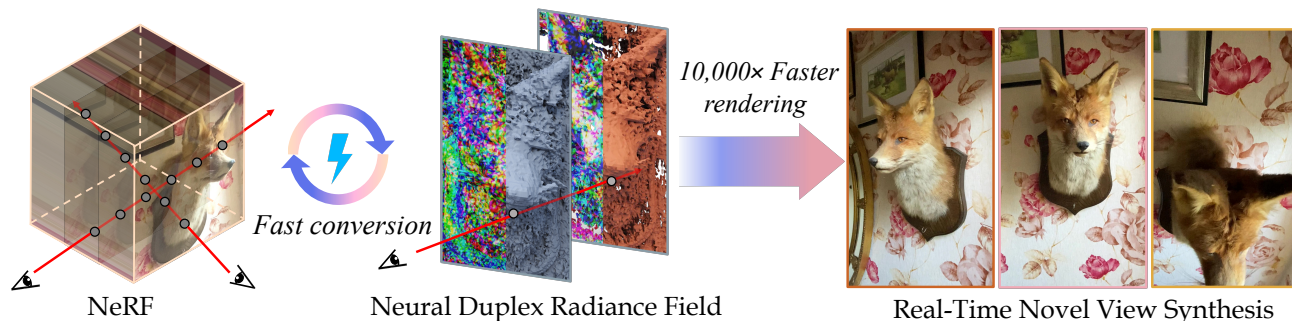


Figure 1. Our framework efficiently learns the neural duplex radiance field from a NeRF model for high-quality real-time view synthesis.

Abstract

Neural radiance fields (NeRFs) enable novel-view synthesis with unprecedented visual quality. However, to render photorealistic images, NeRFs require hundreds of deep multilayer perceptron (MLP) evaluations – for each pixel. This is prohibitively expensive and makes real-time rendering infeasible, even on powerful modern GPUs. In this paper, we propose a novel approach to distill and bake NeRFs into highly efficient mesh-based neural representations that are fully compatible with the massively parallel graphics rendering pipeline. We represent scenes as neural radiance features encoded on a two-layer duplex mesh, which effectively overcomes the inherent inaccuracies in 3D surface reconstruction by learning the aggregated radiance information from a reliable interval of ray-surface intersections. To exploit local geometric relationships of nearby pixels, we leverage screen-space convolutions instead of the MLPs used in NeRFs to achieve high-quality appearance. Finally, the performance of the whole framework is further boosted by a novel multi-view distillation optimization strategy. We demonstrate the effectiveness and superiority of our approach via extensive experiments on a range of standard datasets.

1. Introduction

Reconstructing 3D scenes by a representation that can be rendered from unobserved viewpoints using only a few posed

images has been a long-standing goal in the computer graphics and computer vision communities. Significant progress has recently been achieved by neural radiance fields (NeRFs) [26], which are capable of generating photorealistic novel views and modeling view-dependent effects such as specular reflections. In particular, a radiance field is a volumetric function parameterized by MLPs that estimates density and emitted radiance at sampled 3D locations in a given direction. Differentiable volume rendering then allows the optimization of this function by minimizing the photometric discrepancy between the real observed color and the rendered color.

Despite the unprecedented success and enormous practical potential of NeRF and its various extensions [3, 6, 57], an inescapable problem is the high computational cost of rendering novel views. For instance, even using a powerful modern GPU, NeRF requires about 30 seconds to render a single image with 800×800 pixels, which prevents its use for interactive applications in virtual and augmented reality. On the other hand, the rapid development of NeRF has spawned abundant follow-up works that focus on optimization acceleration [18, 28, 54], generalization [5, 47, 53], and enabled different downstream tasks, including 3D stylization [11, 16, 30], editing [24, 50, 55, 56], or even perception [10, 17]. Thus, a generalized method that can learn and extract a real-time renderable representation given an arbitrary pretrained NeRF-based model, while maintaining high rendering quality, is highly desirable.

The huge computational cost of rendering a NeRF representation mainly comes from two aspects: (1) For each individual pixel, NeRF requires sampling hundreds of loca-

* Corresponding author.

tions along the corresponding ray, querying density and then accumulating radiance using volume rendering; and (2) A large model size is required to represent the geometric details of complex scenes well, so the MLPs used in NeRF architecture are relatively deep and wide, which incurs significant computation for the evaluation of each point sample.

In this paper, we present an approach that achieves high-fidelity real-time view synthesis by addressing these issues. To avoid the dense sampling along each ray, one solution is to generate a geometry proxy from a pretrained NeRF model, e.g., via marching cubes [25]. By exploiting the highly-optimized graphics pipeline, we could almost instantly obtain the sample location for each ray. However, due to inaccuracies in the density field around surfaces, the extracted mesh may not faithfully represent the true underlying geometry and can contain artifacts, as shown in Figure 2. Dense local sampling could alleviate these errors to some degree, but cannot handle missing geometries or occlusions. An alternative approach is to use rasterization for fast neural rendering [1, 37, 38, 43] by directly baking neural features onto the surface of explicit the geometry or point cloud. This is usually accompanied by a deep CNN to translate rasterized features to colors and learn to resolve the existing errors, which is expensive to evaluate and can prevent real-time rendering.

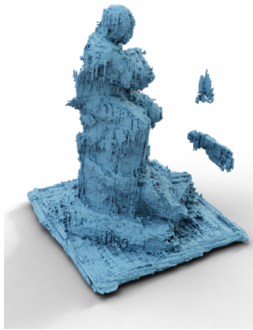


Figure 2. NeRF surface.

To significantly reduce the sampled numbers while efficiently handling the geometric errors, our *first* key technical innovation is to infer the final RGB color according to the radiance of duplex points along the ray. As illustrated in Figure 3, NeRFs represent a continuous density distribution along each ray. While it will generally be difficult to determine the exact location of a surface along the ray, it will be easier to extract a reliable interval that contributes the most to the final prediction by using an under- and an overestimation of the geometry. Motivated by this idea, instead of selecting a specific location for appearance calculation, or performing expensive dense volume rendering in this interval, we represent the scene using learnable features at the two intersection points of a ray with the duplex geometry and use a neural network to learn the color from this aggregated duplex radiance information. Although only two sampled locations are considered, we found this proposed neural duplex radiance field to be robust in compensating for the errors of the geometry proxy even without a deep neural network, while effectively preserving the efficiency of rasterization-based approaches. The NeRF MLP has become the most standard architecture for most neural implicit representations [7, 8, 26, 27, 31, 39]. Yet, with only a few points

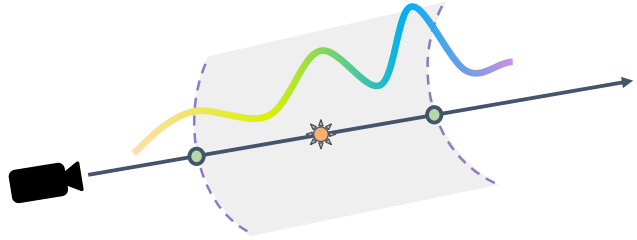


Figure 3. **Motivation.** The continuous density distribution of NeRF (curve above the ray) makes it difficult to identify the accurate location of a surface ($*$) for appearance calculation. We thus seek to extract a reliable interval from the density field (between dashed lines), and learn the duplex radiance combinations to tolerate errors.

considered along the ray, the MLP struggles to constrain the proposed neural duplex radiance field. *Instead*, we use a shallow convolutional network, which can effectively capture the local geometric information of neighboring pixels, and leads to a considerably better rendering quality. *Finally*, we found that directly training the neural duplex radiance field from scratch will lead to noticeable artifacts. We therefore propose a multi-view distillation optimization strategy that enables us to effectively approximate the rendering quality of the original NeRF models. Remarkably, our method improves run-time performance by *10,000* times compared to the original NeRF while maintaining high-quality rendering.

2. Related Work

Our work is in the category of neural rendering. Since the introduction of the seminal NeRF approach [26], the growth of neural rendering papers is exponential [42]. Here, we focus on works closely related to our method along two directions: *neural representations* and *real-time neural rendering*.

Neural Representations. NeRF [26] represents the scene as a continuous field of density and color, and generates images by volume rendering. Built upon NeRF, Mip-NeRF [3] represents the prefiltered radiance field for a continuous space of scales, which enables accuracy improvements for different resolutions. NeRF++ [57] and Mip-NeRF 360 [3] further extend NeRF-like models for unbounded scenes, allowing reconstruction of scenes with far-away backgrounds. NeuS [46] and VolSDF [51] propose to learn a signed distance function representation by volume rendering to better reconstruct 3D objects. This has potential for fast surface rendering, but the visual quality may not reach the photorealism of NeRF. PixelNeRF [53], IBRNet [47] and MVNeRF [5] learn more generalized neural representations to overcome NeRF’s inability to share knowledge between scenes and can achieve high-quality novel-view synthesis given sparse input images. More recently, efforts are made to extremely reduce the training time of NeRF such as Instant-NGP [28], using multiresolution hash encoding for training acceleration, or

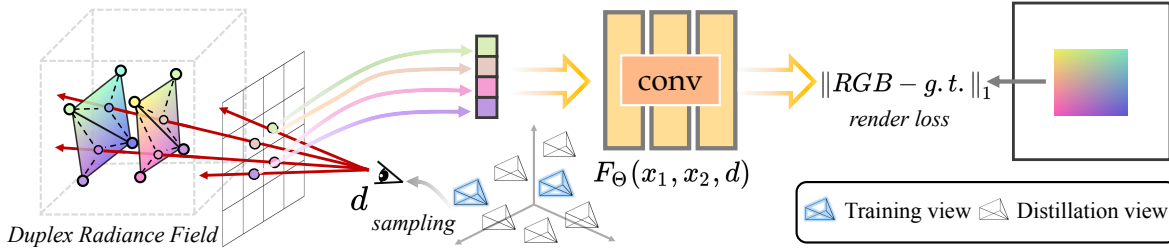


Figure 4. **Architecture of Neural Duplex Radiance Fields.** Our neural representation builds on a two-layer duplex surface extracted from a NeRF model, which approximates the reliable interval covering the true underlying geometry. For each sampled target camera ray with view direction d , we find the duplex ray-surface intersections and derive the radiance features using barycentric interpolation. We use a shallow CNN to aggregate the local geometry and appearance information, and produce view-dependent colors, which effectively ensures rendering efficiency. The output color is directly compared with the pixels of rendered or ground-truth views to optimize the whole representation.

TensorRF [6], using tensor factorization to enable memory-efficient feature storage and so on [18, 41, 54]. Neural light field networks [40] learn to directly map rays to observed radiance, and thus only require a single evaluation per pixel. Recently, multiple ways have been proposed to parameterize camera rays, such as Plücker coordinates [40], Gegenbauer polynomials [12], or local affine embeddings [2]. Albeit with great progress, neural light fields still struggle to achieve view-consistent and high-quality 360-degree rendering due to the lack of an explicit geometry prior.

Real-Time Neural Rendering. There are various directions to accelerate NeRF rendering by reducing the considerable computation needed for inference of large networks and dense sampling along view rays. For example, by deriving sparse intermediate representations from NeRF, such as sparse octrees [45, 52] or sparse voxel grids [14], by dividing the space into many smaller MLPs [34, 35], or by adopting an efficient sampling strategy [15, 21, 29]. Inspired by the rendering equation in computer graphics, FastNeRF [13] proposes to factorize the NeRF and uses caching to achieve efficient rendering. These voxel-based methods basically leverage the caching concept to skip abundant computations in exchange for a larger memory footprint. NeX [49] proposes to combine the MPI representation and neural basis functions to enable real-time view synthesis for forward-facing data. EG3D [4] targets 3D-aware synthesis and supports real-time application by low-resolution rendering followed by upsampling. ENeRF [23] tries to tackle interactive free-viewpoint video by skipping the sampling of empty space. More recently, concurrent work MobileNeRF [9] utilizes optimized sheet meshes as proxy geometry to achieve highly efficient rendering on commodity graphics hardware. In contrast to MobileNeRF [9], our approach neither requires complex multi-stage optimization nor retraining the scene from scratch; instead, our approach can serve as a post-processing step for any existing NeRFs. Another direction explores alternative representations instead of functional NeRFs, such as point clouds [1, 33, 38] or spheres [22].

3. Method

Our main goal is to learn an efficient 3D scene representation that enables high-quality, real-time novel-view synthesis based on a given NeRF model. To this end, we introduce the neural duplex radiance field (NDRF), a novel concept to significantly decrease the number of sampled points along a ray from hundreds to two, while preserving the realistic details of complex scenes well (Section 3.1). Although the query frequency is bounded by using a two-layered mesh, the per-sample deep MLP-based forward pass in NeRF is still burdensome and lacks expressivity, especially for high-resolution rendering scenarios. Thus, we propose a more efficient hybrid radiance representation and shading mechanism (Section 3.2). To further suppress potential artifacts and improve the rendering quality, we leverage multi-view distillation information to guide the optimization procedure (Section 3.3). Finally, we demonstrate the complete rendering procedure and our shader-based implementation that enables real-time cross-platform rendering (Section 3.4).

3.1. Neural Duplex Radiance Fields

Let us start with a quick review of NeRF basics. In a nutshell, NeRF approximates the 5D plenoptic function with a learnable MLP, which maps a spatial location $\mathbf{x} \in \mathbb{R}^3$ and view direction $\mathbf{d} \in \mathbb{S}^2$ to radiance $\mathbf{c} \in \mathbb{R}^3$ and density $\sigma \in \mathbb{R}$. After casting a ray from the camera center through a pixel into the scene, its color is calculated by evaluating a network F_Θ at multiple locations \mathbf{x} along the ray, and aggregating the radiance using volume rendering [26].

A key factor that impacts the rendering efficiency of NeRF is the number of sampled 3D locations for each ray. Since the volumetric representation of NeRF can be converted into a triangle mesh via marching cubes, a natural idea is to directly bake the appearance on a geometric surface, *i.e.*, to only consider the radiance of ray-surface intersections. However, due to the discretization error of marching cubes and the inherent floater artifacts of NeRF [3], the generated mesh

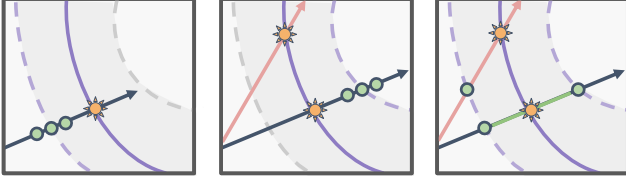


Figure 5. **Duplex radiance field vs local sampling.** \star : real ray-mesh intersection. \bullet : local sampled points around proxy surface. Dashed curve: estimated surface. **Solid curve**: ground-truth surface. **Gray curve**: not currently estimated surface. **Left**: The estimated surface occludes the appearance location on the actual surface. **Middle**: The **red ray** may not hit the estimated surface, resulting in hole artifacts or incorrect intersections. **Right**: Our duplex radiance field instead learns to combine locations on two estimated surfaces to render the correct color, and effectively decreases ray misses.

quality is often inadequate (see Figure 2). Rasterization-based neural rendering methods overcome imperfections in the 3D proxy by correcting errors in screen space using a deep CNN, which inevitably sacrifices rendering efficiency [36, 37]. Since the geometric proxy already conveys the approximate 3D structure, an alternative way would be locally sampling around the surface. Although this can help counter the low-quality mesh, it struggles with missing geometry or visible occlusions, as illustrated in Figure 5.

Representation. It is worth noting that the density distribution of NeRF, which determines the contribution of emitted radiance to the final color of a ray, may contain multiple peaks. This makes approximating the color of a pixel using a single ray-surface intersection a challenging proposition. To significantly reduce the sampled points for each pixel and simultaneously ensure the rendering fidelity, we propose the neural duplex radiance field $F_{\Theta} : \mathbb{R}^8 \mapsto \mathbb{R}^3$ to represent a light ray, where F_{Θ} learns to map a ray segment defined by its two endpoints, \mathbf{x}_1 and \mathbf{x}_2 , to integrated view-dependent colors. More specifically, our goal is to directly learn the color of any pixel from the ray segment defined by the reliable interval of the density distribution. We acquire the endpoints of a ray segment by casting the ray into two proxy mesh surfaces, the inner $\mathcal{M}_i(\mathcal{V}_i, \mathcal{T}_i)$ and outer $\mathcal{M}_o(\mathcal{V}_o, \mathcal{T}_o)$, with vertex positions \mathcal{V}_{\bullet} and triangle faces \mathcal{T}_{\bullet} . These meshes are extracted from different level sets of NeRF’s density field. More formally, the ray color \mathbf{c} could be calculated by

$$\mathbf{c} = F_{\Theta}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{d}) \quad \text{for } \mathbf{x}_1 \in \mathcal{T}_o \text{ and } \mathbf{x}_2 \in \mathcal{T}_i, \quad (1)$$

where \mathbf{x}_1 and \mathbf{x}_2 are the depth-sorted intersections of a given ray with meshes \mathcal{M}_i and \mathcal{M}_o , and \mathbf{d} denotes the view direction. We still append the view direction considering not all rays will hit both surface proxies. The explicit triangle mesh then will allow us to efficiently acquire the ray-surface intersection, benefiting from the full parallelism of the graphics rasterization pipeline.

In most implicit neural representations, different properties like SDF, radiance, or pixel colors are encoded by a

multilayer perceptron due to its smoothness and compactness properties. However, densely evaluating a deep MLP is not a wise choice for real-time applications considering the high computational cost of inference. Since our radiance field is defined via two geometric surfaces, our encoding is defined by directly attaching learnable features $\mathbf{f}_k \in \mathbb{R}^N$ to each mesh vertex \mathbf{v}_k . To obtain features $\mathbf{f}(\mathbf{x})$ for any point \mathbf{x} on a triangle, we interpolate features within triangles using barycentric coordinates via hardware rasterization. Thus, the rendering equation can be re-written as

$$\mathbf{c} = F_{\Theta}(\mathbf{f}(\mathbf{x}_1), \mathbf{f}(\mathbf{x}_2), \mathbf{d}) \quad \text{with } \mathbf{x}_1 \in \mathcal{T}_o, \mathbf{x}_2 \in \mathcal{T}_i, \quad (2)$$

where Θ parameterizes a shallow network to learn the aggregation of radiance features given each ray.

3.2. Convolutional Shading

Many NeRF-based methods use a view direction conditioned MLP to learn view-dependent appearance effects. The dense sampling along rays provides a feasible way to constrain the learning of the radiance fields through shared 3D locations, which ultimately enables high-quality and coherent rendering across different views. This implicit constraint, however, does not apply to our case. To highly optimize run-time efficiency, neural duplex radiance only considers two points for each cast ray, which significantly reduces the capability of capturing internal correlation in training, often resulting in spatial artifacts and the degradation of rendering quality.

We resolve this issue by a convolutional shading network that converts duplex features and view directions to per-pixel colors in image space using a small receptive field. More specifically, an independent pixel will be rendered by considering the neighboring local geometry positions \mathbf{x}_i whose projection $\mathbf{K}(\mathbf{R}\mathbf{x}_i + \mathbf{t})$ with the nearest z -coordinate lies in the local window around this pixel location, where camera pose $\mathbf{R} \in \text{SO}(3)$, $\mathbf{t} \in \mathbb{R}^3$ and \mathbf{K} is the intrinsic matrix. This approach effectively increases the correlation of samples of a 3D scene. To ensure high inference speed, we only use 2~3 convolutional layers and limit convolution kernels to a small 2×2 window for optimal visual fidelity and view consistency. Please refer to the supplemental material for more detailed network architecture and parameters.

3.3. Multi-View Distillation

We observed that neural duplex radiance fields trained from scratch fail to match the high fidelity of the original NeRFs and produce many spatial high-frequency artifacts. To overcome this problem, we take inspiration from KiloNeRF [35] and use a pretrained NeRF as a teacher model to guide the optimization of our neural duplex radiance fields. KiloNeRF directly minimizes the differences of density and radiance in 3D volume space between teacher and student models without any rendering. Considering different parameterizations

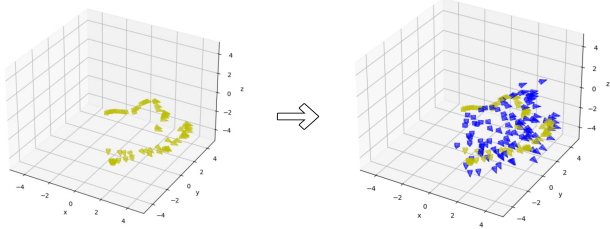


Figure 6. **Visualizations of distillation view generation.** Yellow cameras: training views. Blue cameras: sampled distillation views.

of NeRF and our models, we instead use rendered multi-view images to distill knowledge from the teacher model.

One important question is how to generate effective distillation views for rendering based on the training dataset. For object-level datasets, camera views tend to lie on a hemisphere that can be estimated and used for sampling new distillation viewpoints similar with R2L [44]. For real-world captures, the camera poses may only occupy a small subspace. To produce meaningful interpolated viewpoints under both settings, we first align the view directions of all cameras with the origin and convert the camera poses into Spherical coordinates (r, θ, φ) relative to the origin. For each distillation view, we then randomly sample a radius $r_s \sim U(r^{\min}, r^{\max})$, angles $\theta_s \sim U(\theta^{\min}, \theta^{\max})$ and $\varphi_s \sim U(\varphi^{\min}, \varphi^{\max})$ according to the ranges of spherical coordinates. After converting from spherical to Cartesian coordinates, we obtain many suitable interpolated poses following the original viewpoint distribution. We define the number of distillation views to be 1,000 in all settings. Our model will be trained on the distilled views first and then fine-tuned using the training images to produce sharper appearance.

3.4. Real-Time Rendering

A fully trained neural duplex radiance field comprises two meshes with per-vertex features as well as the shallow shading CNN. While we train the model using PyTorch [32], we implemented a real-time renderer in WebGL using GLSL shaders for cross-platform compatibility. To synthesize a novel view, we use hardware rasterization to efficiently render two feature buffers at the output image resolution. These buffers are sent into the CNN together with a per-pixel view direction to aggregate the local radiance features and produce view-dependent RGB colors. For efficiency, we implement the CNN in two rendering passes, one pass for each layer. Thanks to the compact size of the convolution kernels, each convolution layer can be executed in a pixel-parallel way, which is highly efficient, even without CUDA.

4. Experiments

4.1. Setup

Datasets. We evaluate and compare performance on the NeRF-Synthetic [26] benchmark dataset, which contains

eight objects with 360° views at 800×800 resolution. Further, to comprehensively compare the quality and to test the upper bound of speed, we conduct experiments on the megapixel-level dataset of Tanks&Temples [20], which contains 1920×1080 frames. We also leverage additional real-world datasets to perform ablation studies. All the training and testing splits follow the published literature.

Baselines. We compare with two representative efficient novel-view synthesis approaches: KiloNeRF [35] and SNeRG [14]. For quantitative comparisons, we directly use the published numerical results. For qualitative comparisons, we retrain both methods on different datasets following original configurations, since not all checkpoints are released.

Evaluation Metrics. We measure rendering quality by comparing predicted novel views and ground-truth images using standard metrics: peak signal-to-noise ratio (PSNR), structural similarity index (SSIM) [48], and learned perceptual image patch similarity (LPIPS) [58]. Since our goal is highly efficient rendering, we also report framerates using frames per second (FPS). We also compare the GPU memory consumption in GB for WebGL applications, to quantify the trade-off between speed, quality, and memory. Finally, we compare the training speed by recording the total optimization time in hours. Our test system consists of an NVIDIA RTX 2080 Ti GPU, an i7-9700K CPU, and 32 GB RAM.

Implementation. We implement the CUDA version with 20-dimensional learnable radiance features and a three-layer convolutional network. To ensure high efficiency and match the shader implementation, in the WebGL version, we decrease the learnable feature dimension attached to the geometry surface from 20 to 8 and remove the final convolution layer. This results in some performance degradation on both datasets compared with our CUDA version, but brings remarkable FPS improvement and cross-platform properties, since the renderer implementation no longer requires CUDA. More details regarding the architecture design and parameters setting can be found in the supplementary material.

Optimization Details. We build our method on top of PyTorch [32] and use the Adam optimizer [19] with $\beta_1 = 0.9$ and $\beta_2 = 0.999$. Our method converges in 200,000 iterations, where the first half is optimized on the distillation views and the second half on the original training images. In each iteration, we randomly sample one view, cast all rays into the geometry to fetch the duplex radiance features, and generate complete screen-space buffers (or local patches if training memory is insufficient) for the convolutional shading. The learning rate is set to 10^{-3} initially, with exponential decay to 10^{-5} thereafter. We chose TensorRF [6] as the teacher NeRF model considering its high efficiency of reconstruction and synthesis of distillation views. To convert the density field into the duplex mesh geometry, the two thresholds are set empirically to 10^{-4} and 10^{-2} , respectively.

Table 1. **Quantitative results on Synthetic-NeRF [26] and Tanks&Temples [20] datasets.** The best result is highlighted in **bold**. Ours (*laptop*): Run on an M1 MacBook Pro with 16 GB RAM. Others run on a desktop. N+M: N is the training time and M is the baking time.

Dataset	Method	Backend	FPS \uparrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Training (hours)	Training Specs
Synthetic-NeRF 800×800	KiloNeRF [35]	CUDA	33.47	31.00	0.950	0.030	~35+14	1×V100
	SNeRG [14]	WebGL	89.43	30.38	0.950	0.050	~30+5	8×A100
	Ours	CUDA	59.71	32.14	0.957	0.030	6.61	1×3090
	Ours	WebGL	282.67	30.43	0.945	0.049	1.01	1×3090
	Ours (<i>laptop</i>)	WebGL	62.27	30.43	0.945	0.049	1.01	1×3090
Tanks&Temples 1920×1080	KiloNeRF	CUDA	16.30	28.41	0.910	0.090	~24+24	1×A100
	SNeRG	WebGL	55.48	24.95	0.882	0.184	~30+5	8×A100
	Ours	CUDA	27.49	28.12	0.915	0.089	9.02	1×A100
	Ours	WebGL	186.94	27.08	0.895	0.130	1.84	1×A100
	Ours (<i>laptop</i>)	WebGL	33.64	27.08	0.895	0.130	1.84	1×A100

Table 2. **Memory–Quality–Speed trade-off.** Peak GPU memory consumption (GB), rendering quality, run-time performance comparisons on Tanks&Temples [20] dataset for WebGL applications.

Method	Memory (GB) \downarrow	PSNR \uparrow	FPS \uparrow
SNeRG	4.19	24.95	55.48
Ours	0.99	27.08	186.94
Ours (<i>laptop</i>)	0.94	27.08	33.64

Table 3. **Quantitative ablation study of proposed components.** Our full model performs significantly better than all variations.

Duplex	Conv	Distill	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
	✓	✓	26.95	0.884	0.12
✓		✓	30.03	0.942	0.05
✓	✓		29.66	0.932	0.06
✓	✓	✓	32.14	0.957	0.03

4.2. Results

Quantitative Comparison. We provide extensive quantitative comparisons in Table 1. Compared with CUDA-based KiloNeRF [35], our CUDA version achieves significantly better rendering performance in terms of PSNR (more than 1 dB) and SSIM on the Synthetic-NeRF dataset [26] with much higher run-time FPS and lower training time. The PSNR of our CUDA-based model on Tanks&Temples [20] is slightly lower than KiloNeRF, even though the novel views of KiloNeRF contain severe artifacts, as shown in Figure 7. This may be caused by the characteristics of PSNR, which only measures pixel-level differences and ignores the structural similarity. In this situation, SSIM would be a better metric to indicate the performance gap between our method and KiloNeRF (0.915 vs. 0.910).

We also compare with WebGL-based SNeRG [14]. Even at 1920×1080 resolution, our WebGL-based model achieves 180+ FPS on a desktop-level GPU, and 30+ FPS on a laptop

using Chrome, with better rendering fidelity than SNeRG. In addition, our method requires about two orders of magnitude less computation for training compared to SNeRG. We also show the trade-off comparison regarding memory consumption, image quality and rendering speed in Table 2 under high-resolution setting, which comprehensively demonstrates the superiority of our method.

Qualitative Comparison. To further verify our method’s benefit in terms of rendering quality over existing fast rendering NeRF techniques, we show qualitative comparisons in Figure 7. As we can see, both baselines, KiloNeRF and SNeRG, cannot render the high-quality details of the scenes, which are possibly bounded by the voxel resolution; however, simply increasing the voxel resolution will directly lead to cubic growth of memory cost. Besides, some texture artifacts and wrong view-dependent specular highlights also exist in the baselines, which are visually inferior to our method. More importantly, in the Tanks&Temples dataset, we could observe vibrant blocking degradation from the results of KiloNeRF. The potential reason is KiloNeRF has limited generalization capability to out-of-distribution (OOD) views since there are no correlations between space-isolated MLPs. By contrast, our method can render consistent and high-fidelity views with good view generalization ability.

4.3. Ablation Study

Duplex Radiance. To verify the effectiveness of our proposed neural duplex radiance, we compare with a variation that only considers a single ray-surface intersection. More specifically, we only adopt the mesh extracted from NeRF model with threshold 10^{-4} to ensure geometry proxy could cover all pixels and avoid hole artifacts. As shown in Table 3 and Figure 8, rendering without duplex radiance dramatically downgrades the quality of novel-view synthesis since the coarse geometry extracted from NeRF could not provide reasonable shading locations. More results about using single

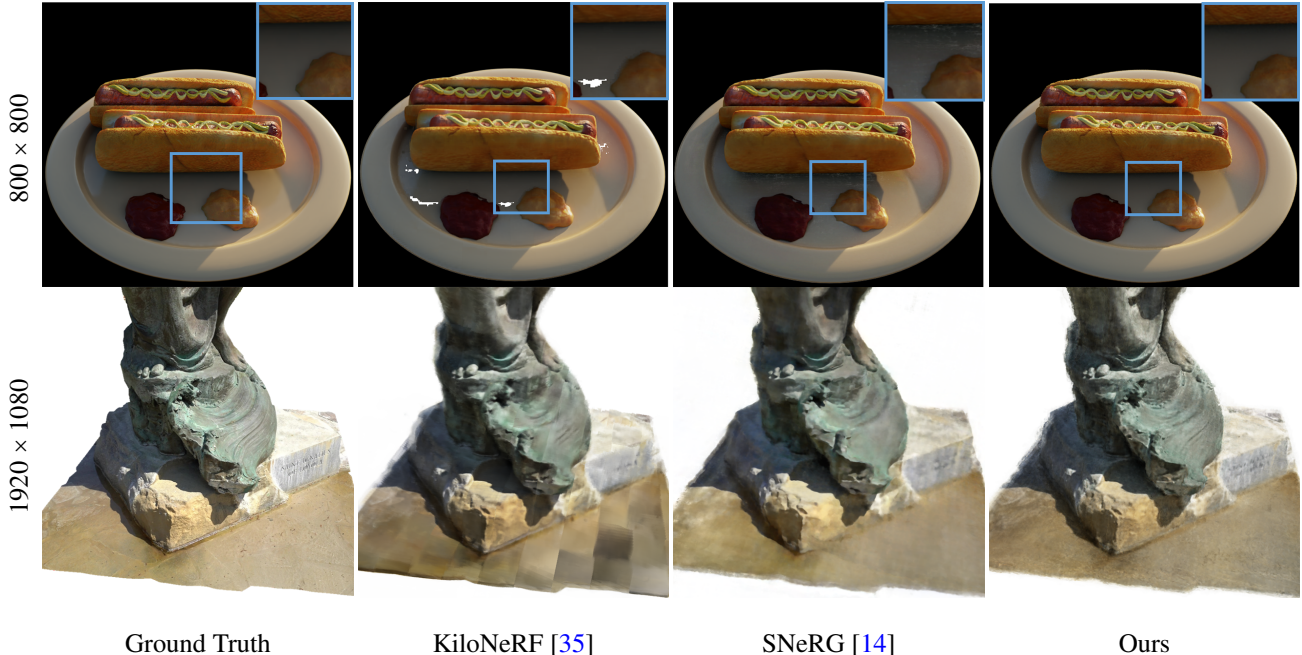


Figure 7. **Qualitative comparisons with various baselines on Synthetic-NeRF (top) and Tanks&Temples (bottom) datasets.** We show the novel views synthesized by KiloNeRF, SNeRG and ours. Our approach effectively speeds up the inference efficiency while maintaining the high quality of the rendered frames.

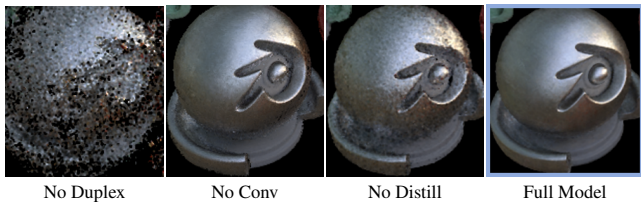


Figure 8. **Qualitative results for each ablation study.** Our full model achieves the best rendering quality among all variations.

radiance feature but with different thresholds can be found in the supplementary materials.

Convolution vs MLP. In contrast to most NeRF methods, which integrate radiance only along a single ray, we aggregate all the radiance information of the local geometry surface via convolution kernels. Figure 8 and Table 3 demonstrate that the convolutional shading network can render high-fidelity novel views more robustly, and effectively overcomes the under-constrained problem while only considering two sample locations for each ray.

No Distillation. In this ablation study, we remove all the distillation views and only train the model based on known views. Our results in Figure 8 show that optimization without distillation views will lead to high-frequency artifacts in the geometry, which the spatial convolutional kernels are unable to effectively handle as well.

Influence of Network Size. We have noticed that decreasing the radiance feature dimension and network size will result in some quality degeneration in Table 1. Hence, one

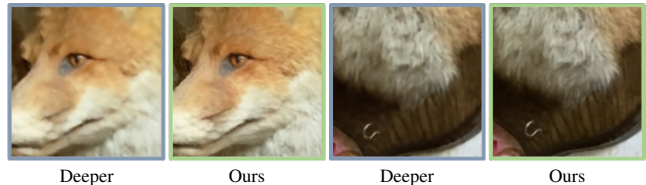


Figure 9. **Ablation study on the network size.** Increasing the depth and receptive field of the convolutional shading network does not lead to an improvement on rendering quality in a consistent manner. Zoom in for more details.

interesting question is whether performance can be boosted consistently by increasing the network parameters and receptive fields. We conduct this experiment on the Fox dataset. More specifically, we double the network layers and boost the original used convolution layers with larger receptive fields. As shown in Figure 9, the deeper network does not bring more quality gains. By contrast, the synthesized novel views become more blurry than our method and lose many textural details, which may be caused by the learning difficulty of a deep network. Leveraging more advanced architectures or loss functions, such as adversarial objectives, could potentially resolve this issue.

Generalization to out-of-distribution views. NeRF has strong capabilities to interpolate novel views, but its rendering quality may drop for some novel views which don't follow the original distribution. This issue will be more serious for distillation-based methods that try to mimic the teacher behavior. We show out-of-distribution (OOD) syn-



Figure 10. **Comparison of out-of-distribution (OOD) view synthesis (No ground-truth).** Our method could generate high-quality and view-dependent frames with sharper details for OOD views.

Table 4. **Quantitative comparisons with the teacher model** on the Tanks&Temples [20] dataset (1920×1080).

	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FPS \uparrow
Original NeRF [26]	28.32	0.900	0.11	0.01
Teacher TensorRF [6]	28.48	0.918	0.12	0.10
Ours w/ WebGL	27.08	0.895	0.13	186.94
Ours w/ CUDA	28.12	0.915	0.09	27.49

thesis results in Figure 10. Both KiloNeRF [35] and SNeRG [14] do not generalize well to OOD views and produce inconsistent and blurry artifacts, while our method synthesizes sharper and cleaner appearance with appropriate lighting.

Comparison with Teacher Model. Since our method involves multi-view distillation, we also measure the quantitative differences between our model and the teacher model [6] in Table 4. Although the numerical results of PSNR and SSIM are slightly inferior to the teacher model, the perceptual metric LPIPS, which better correlates with human perception, surprisingly attains better performance benefiting from the better generalization ability of neural duplex radiance field. On the other hand, the teacher model TensorRF [6] has already optimized the run-time performance by skipping the importance sampling step of original NeRF model and tensor decomposition, but the speed is still limited for real-world interactive applications. In contrast, our method achieves 30+ FPS while maintaining comparable rendering quality. Significantly, our WebGL version reaches 10,000× speed up over the original NeRF model.

Comparison with MobileNeRF [9]. Recently, a concurrent work MobileNeRF [9] also explores the rasterization pipeline to accelerate the rendering speed of neural radiance fields. Nonetheless, their framework involves optimizations of both geometry and appearance leading to a high computational cost (21+ hours using 8×V100). By contrast, our method can convert a pretrained NeRF in a few hours using a single GPU. We provide some qualitative comparisons in Figure 11 as well. Our approach produces more appealing results and better view-dependent appearance.

5. Conclusion

In this paper, we proposed a novel approach for learning a neural duplex radiance field from a pretrained NeRF for real-

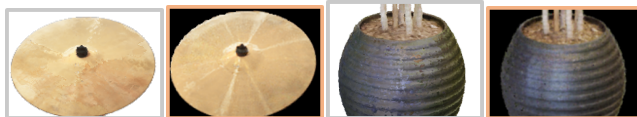


Figure 11. **Qualitative comparisons with MobileNeRF.** Our method produces smoother appearance.

time novel-view synthesis. We employed a convolutional shading network to improve the rendering quality and also proposed a multi-view distillation strategy for better optimization. Our method provides significant improvements on run-time, bettering or preserving the rendering quality, and at the same time, operating on a much lower computational cost compared to existing efficient NeRF methods.

Limitations. It is challenging to extract useful duplex meshes from NeRFs to represent transparent or semi-transparent scenes and our method is no different. Integrating order-independent transparency (OIT) with the NeRF framework may be a potential solution to speed up its rendering. In addition, the quality of the geometry proxy will directly influence the learning of neural representations, thus we would also like to explore scene-specific thresholds rather than empirical defined values while generating meshes, for instance, based on signed distance functions or on the statistics of the density values, for more robust proxy extraction. Currently, our method focuses on object-level and bounded-scene datasets. Efficiently rendering real-world unbounded scenes is a promising future direction. Finally, our experiments demonstrate that sufficient learnable parameters are crucial to achieve better rendering fidelity. Hence, to counter the performance degradation in the WebGL version, boosting the running efficiency of larger CNN in the fragment shader, or leveraging other acceleration frameworks, could also be considered in future work.

Acknowledgements. We thank the anonymous reviewers for their constructive comments. We also appreciate helpful discussions with Feng Liu, Chakravarty R. Alla Chaitanya, Simon Green, Daniel Maskit, Aayush Bansal, Zhiqin Chen, Vasu Agrawal, Hao Tang, Michael Zollhoefer, Huan Wang, Ayush Saraf and Zhaoyang Lv. This work was partially supported by a GRF grant (Project No. CityU 11216122) from the Research Grants Council (RGC) of Hong Kong.

References

- [1] Kara-Ali Aliev, Artem Sevastopolsky, Maria Kolos, Dmitry Ulyanov, and Victor Lempitsky. Neural point-based graphics. In *ECCV*, pages 696–712, 2020.
- [2] Benjamin Attal, Jia-Bin Huang, Michael Zollhöfer, Johannes Kopf, and Changil Kim. Learning neural light fields with ray-space embedding. In *CVPR*, pages 19819–19829, 2022.
- [3] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-NeRF 360: Unbounded anti-aliased neural radiance fields. In *CVPR*, pages 5470–5479, 2022.
- [4] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *CVPR*, pages 16123–16133, 2022.
- [5] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. MVSNerF: Fast generalizable radiance field reconstruction from multi-view stereo. In *CVPR*, pages 14124–14133, 2021.
- [6] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. TensoRF: Tensorial radiance fields. In *ECCV*, 2022.
- [7] Yinbo Chen, Sifei Liu, and Xiaolong Wang. Learning continuous image representation with local implicit image function. In *CVPR*, pages 8628–8638, 2021.
- [8] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *CVPR*, pages 5939–5948, 2019.
- [9] 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.
- [10] Nianchen Deng, Zhenyi He, Jiannan Ye, Budmonde Duinkharjav, Praneeth Chakravarthula, Xubo Yang, and Qi Sun. FoV-NeRF: Foveated neural radiance fields for virtual reality. *TVCG*, 28(11):3854–3864, 2022.
- [11] Zhiwen Fan, Yifan Jiang, Peihao Wang, Xinyu Gong, DeJia Xu, and Zhangyang Wang. Unified implicit neural stylization. In *ECCV*, 2022.
- [12] Brandon Yushan Feng and Amitabh Varshney. Signet: Efficient neural representation for light fields. In *ICCV*, pages 14224–14233, 2021.
- [13] Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. FastNeRF: High-fidelity neural rendering at 200fps. In *CVPR*, pages 14346–14355, 2021.
- [14] Peter Hedman, Pratul P Srinivasan, Ben Mildenhall, Jonathan T Barron, and Paul Debevec. Baking neural radiance fields for real-time view synthesis. In *CVPR*, pages 5875–5884, 2021.
- [15] Tao Hu, Shu Liu, Yilun Chen, Tiancheng Shen, and Jiaya Jia. EfficientNeRF – efficient neural radiance fields. In *CVPR*, pages 12902–12911, 2022.
- [16] Yi-Hua Huang, Yue He, Yu-Jie Yuan, Yu-Kun Lai, and Lin Gao. StylizedNeRF: Consistent 3D scene stylization as stylized NeRF via 2D–3D mutual learning. In *CVPR*, 2022.
- [17] Yoonwoo Jeong, Seungjoo Shin, Junha Lee, Chris Choy, Anima Anandkumar, Minsu Cho, and Jaesik Park. PerFception: Perception using radiance fields. In *NeurIPS Datasets and Benchmarks Track*, 2022.
- [18] Animesh Karnewar, Tobias Ritschel, Oliver Wang, and Niloy Mitra. ReLU fields: The little non-linearity that could. In *SIGGRAPH Conference Proceedings*, pages 27:1–9, 2022.
- [19] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [20] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Trans. Graph.*, 36(4):1–13, 2017.
- [21] Andreas Kurz, Thomas Neff, Zhaoyang Lv, Michael Zollhöfer, and Markus Steinberger. AdaNeRF: Adaptive sampling for real-time rendering of neural radiance fields. In *ECCV*, 2022.
- [22] Christoph Lassner and Michael Zollhofer. Pulsar: Efficient sphere-based neural rendering. In *CVPR*, pages 1440–1449, 2021.
- [23] Haotong Lin, Sida Peng, Zhen Xu, Yunzhi Yan, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Efficient neural radiance fields for interactive free-viewpoint video. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–9, 2022.
- [24] Steven Liu, Xiuming Zhang, Zhoutong Zhang, Richard Zhang, Jun-Yan Zhu, and Bryan Russell. Editing conditional radiance fields. In *ICCV*, pages 5773–5783, 2021.
- [25] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3D surface construction algorithm. *SIGGRAPH Computer Graphics*, 21(4):163–169, 1987.
- [26] 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.
- [27] 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. *Communications of the ACM*, 65(1):99–106, 2021.
- [28] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–15, 2022.
- [29] Thomas Neff, Pascal Stadlbauer, Mathias Parger, Andreas Kurz, Joerg H Mueller, Chakravarty R Alla Chaitanya, Anton Kaplanyan, and Markus Steinberger. DONeRF: Towards real-time rendering of compact neural radiance fields using depth oracle networks. *Computer Graphics Forum*, 40(4):45–59, 2021.
- [30] Thu Nguyen-Phuoc, Feng Liu, and Lei Xiao. SNeRF: Stylized neural implicit representations for 3D scenes. *ACM Trans. Graph.*, 41(4):142:1–11, 2022.
- [31] Michael Oechsle, Lars Mescheder, Michael Niemeyer, Thilo Strauss, and Andreas Geiger. Texture fields: Learning texture representations in function space. In *ICCV*, pages 4531–4540, 2019.
- [32] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu

- Fang, Junjie Bai, and Soumith Chintala. PyTorch: An imperative style, high-performance deep learning library. In *NeurIPS*, pages 8024–8035, 2019.
- [33] Ruslan Rakhimov, Andrei-Timotei Ardelean, Victor Lempitsky, and Evgeny Burnaev. NPBG++: Accelerating neural point-based graphics. In *CVPR*, pages 15948–15958, 2022.
- [34] Daniel Rebain, Wei Jiang, Soroosh Yazdani, Ke Li, Kwang Moo Yi, and Andrea Tagliasacchi. DeRF: Decomposed radiance fields. In *CVPR*, pages 14153–14161, 2021.
- [35] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. KiloNeRF: Speeding up neural radiance fields with thousands of tiny MLPs. In *ICCV*, pages 14335–14345, 2021.
- [36] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *ECCV*, 2020.
- [37] Gernot Riegler and Vladlen Koltun. Stable view synthesis. In *CVPR*, pages 12216–12225, 2021.
- [38] Darius Rückert, Linus Franke, and Marc Stamminger. ADOP: Approximate differentiable one-pixel point rendering. *ACM Trans. Graph.*, 41(4):99:1–14, 2022.
- [39] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3D-structure-aware neural scene representations. In *NeurIPS*, 2019.
- [40] Vincent Sitzmann, Semon Rezhikov, William T. Freeman, Joshua B. Tenenbaum, and Fredo Durand. Light field networks: Neural scene representations with single-evaluation rendering. In *NeurIPS*, 2021.
- [41] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *CVPR*, 2022.
- [42] Ayush Tewari, Justus Thies, Ben Mildenhall, Pratul Srinivasan, Edgar Tretschk, Yifan Wang, Christoph Lassner, Vincent Sitzmann, Ricardo Martin-Brualla, Stephen Lombardi, Tomas Simon, Christian Theobalt, Matthias Niessner, Jonathan T. Barron, Gordon Wetzstein, Michael Zollhöfer, and Vladislav Golyanik. Advances in neural rendering. *Computer Graphics Forum*, 41(2):703–735, 2022.
- [43] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *ACM Trans. Graph.*, 38(4):66:1–12, 2019.
- [44] Huan Wang, Jian Ren, Zeng Huang, Kyle Olszewski, Menglei Chai, Yun Fu, and Sergey Tulyakov. R2L: Distilling neural radiance field to neural light field for efficient novel view synthesis. In *ECCV*, pages 612–629, 2022.
- [45] Liao Wang, Jiakai Zhang, Xinhang Liu, Fuqiang Zhao, Yan-shun Zhang, Yingliang Zhang, Minye Wu, Jingyi Yu, and Lan Xu. Fourier PlenOctrees for dynamic radiance field rendering in real-time. In *CVPR*, pages 13524–13534, 2022.
- [46] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. In *NeurIPS*, 2021.
- [47] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P. Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. IBRNet: Learning multi-view image-based rendering. In *CVPR*, pages 4690–4699, 2021.
- [48] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Trans. Image Process.*, 13(4): 600–612, 2004.
- [49] Suttisak Wizadwongsa, Pakkapon Phongthawee, Jiraphon Yenphraphai, and Supasorn Suwajanakorn. NeX: Real-time view synthesis with neural basis expansion. In *CVPR*, pages 8534–8543, 2021.
- [50] Bangbang Yang, Yinda Zhang, Yinghao Xu, Yijin Li, Han Zhou, Hujun Bao, Guofeng Zhang, and Zhaopeng Cui. Learning object-compositional neural radiance field for editable scene rendering. In *ICCV*, pages 13779–13788, 2021.
- [51] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. In *NeurIPS*, pages 4805–4815, 2021.
- [52] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. PlenOctrees for real-time rendering of neural radiance fields. In *ICCV*, 2021.
- [53] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. PixelNeRF: Neural radiance fields from one or few images. In *CVPR*, pages 4578–4587, 2021.
- [54] Alex Yu, Sara Fridovich-Keil, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022.
- [55] Yu-Jie Yuan, Yang-Tian Sun, Yu-Kun Lai, Yuewen Ma, Rongfei Jia, and Lin Gao. NeRF-editing: geometry editing of neural radiance fields. In *CVPR*, pages 18353–18364, 2022.
- [56] Jingbo Zhang, Xiaoyu Li, Ziyu Wan, Can Wang, and Jing Liao. FDNerF: Few-shot dynamic neural radiance fields for face reconstruction and expression editing. In *SIGGRAPH Asia Conference Papers*, pages 12:1–9, 2022.
- [57] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. NeRF++: Analyzing and improving neural radiance fields. arXiv:2010.07492, 2020.
- [58] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, pages 586–595, 2018.

Learning Neural Duplex Radiance Fields for Real-Time View Synthesis

Supplementary Material

Ziyu Wan¹ Christian Richardt² Aljaž Božič² Chao Li² Vijay Rengarajan²
Seonghyeon Nam² Xiaoyu Xiang² Tuotuo Li² Bo Zhu² Rakesh Ranjan² Jing Liao^{1*}

¹City University of Hong Kong ²Meta Reality Labs

1. Overview

In this supplemental material, additional implementation details and experimental results are provided, including:

- More details about network architecture (Section 2);
- More details about shader implementation (Section 3);
- More ablation studies about threshold selection, the teacher model and the number of mesh layers. We also provide another perspective to visualize and understand the proposed method (Section 4);
- Video demo. Please refer to our project page.

2. Network Architecture

We present the details of the two versions of our shading CNN in Table 1. To generate high-frequency textures and view-dependent effects, we also leverage the positional encoding, which maps view direction and intersection locations into a higher-dimensional space. For the CUDA version, both the position and view encoding dimensions are set to 10. For the WebGL version, we set the dimension of view direction encoding to 5 and the others to 0. To ensure the CNN efficiency of the CUDA version, we use depthwise separable convolution [2] to avoid expensive computations. We also tried to implement depthwise separable convolution of WebGL version into the fragment shader but didn't attain effective speedup.

Table 1. Detailed architecture of convolutional shading network.

Version	Layer	Kernel size / stride	Channels	Non-Linearity
CUDA	2DConv	$3 \times 3 / (1, 1)$	256 \rightarrow 256	ReLU
	2DConv	$3 \times 3 / (1, 1)$	256 \rightarrow 256	ReLU
	2DConv	$1 \times 1 / (1, 1)$	256 \rightarrow 3	Sigmoid
WebGL	2DConv	$2 \times 2 / (1, 1)$	55 \rightarrow 32	ReLU
	2DConv	$2 \times 2 / (1, 1)$	32 \rightarrow 3	Sigmoid

3. Shader Implementation

Since the learnable features are attached on the mesh vertices, we directly generate the screen-space feature buffer using traditional hardware rasterization. More specifically, we render 4 RGBA buffers for each mesh: 2 for 8-channel features, 1 for ray-surface intersection positions, and 1 for view directions. To implement the convolution operator in the same framework, after optimizing the whole parameters using PyTorch, we import the 2-layer convolution weights into 2 RGBA `GL_TEXTURE_2D` with shape `out_channel × in_channel` thanks to the 2×2 spatial kernel. We employ two passes in total to generate view-dependent RGB frames, one pass for each convolutional layer. In each pass, we only consider the

* Corresponding author.

Table 2. **Quantitative comparisons on Synthetic-NeRF [3] Ficus data with different thresholds.** Our approach significantly outperforms single-surface based baselines.

	10^{-4}	5×10^{-4}	10^{-3}	5×10^{-3}	10^{-2}	Ours
PSNR \uparrow	25.50	26.91	27.52	28.01	27.05	32.67
SSIM \uparrow	0.921	0.936	0.941	0.944	0.935	0.975
LPIPS \downarrow	0.088	0.073	0.068	0.063	0.069	0.024

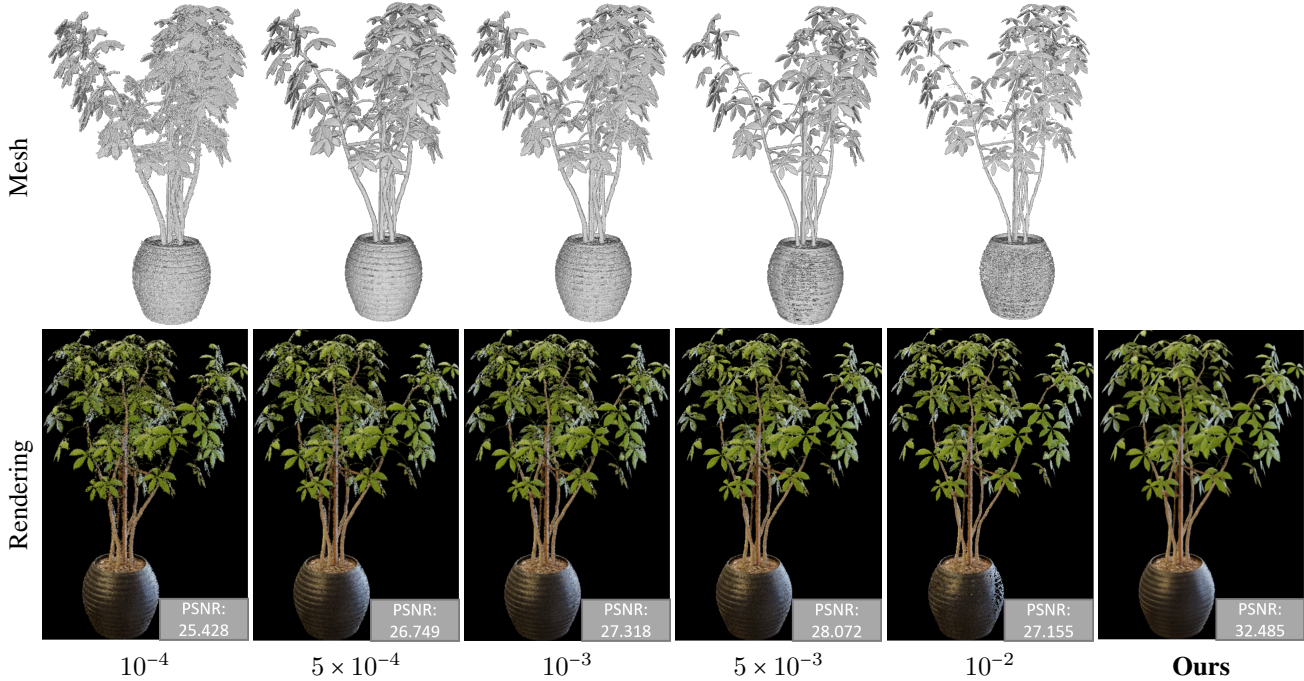


Figure 1. **Qualitative comparisons on Synthetic-NeRF [3] Ficus data with different thresholds.** The corresponding mesh is also visualized in the top row. We observe that it is difficult to employ a single marching cubes surface to faithfully represent a NeRF.

radiance information of a local receptive field, which can be efficiently queried through `texelFetch`. To minimize the memory footprint, we implement `sin/cos` positional encoding in the fragment shader of the first convolution layer rather than in the rasterization step. We additionally tried to put the forward propagation of all CNN layers in one fragment shader, but found the efficiency to be poor. There are also some well-known existing web-based deep learning frameworks like `TF.js`, but we found they could not efficiently support the high-resolution synthesis.

4. Additional Ablation Studies

Results of Single Surface using Different Thresholds. In the main paper, we have conducted ablation studies on the effectiveness of duplex radiance fields by using a single extracted mesh with a threshold 10^{-4} . Here another interesting question is, what will be the best performance if we adjust the threshold to get a better surface? Will this surpass the performance of neural duplex radiance fields? To answer these questions, we conducted both quantitative and qualitative experiments on NeRF-Synthetic Ficus data [3]. As shown in Table 2 and Figure 1, carefully fine-tuning the threshold will lead to some performance improvements, but the rendering quality is still not acceptable due to the inaccurate geometry. In contrast, through learning the aggregation of radiance features on two different coarse geometry surfaces, our method achieves significantly better novel-view synthesis quality than all these variations. Please note the isolated parts of extracted mesh will be automatically removed with respect to their diameters.

Results using Other Teacher Models. To prove the generalization ability of our method, we also conduct experiments on learning neural duplex radiance fields from other NeRF models. More specifically, we train Instant-NGP [4] on the

Table 3. **Quantitative comparisons using different teacher models.** On Synthetic-NeRF [3] Chair data, we show our method have generalization capabilities to different NeRF models.

	TensorRF [1]	Ours-TensorRF	Instant-NGP [4]	Ours-Instant-NGP
PSNR \uparrow	34.73	34.08	34.30	34.17
SSIM \uparrow	0.981	0.983	0.979	0.982
LPIPS \downarrow	0.013	0.013	0.010	0.011



Figure 2. **Qualitative comparisons on Synthetic-NeRF [3] Chair data using different teacher models.**

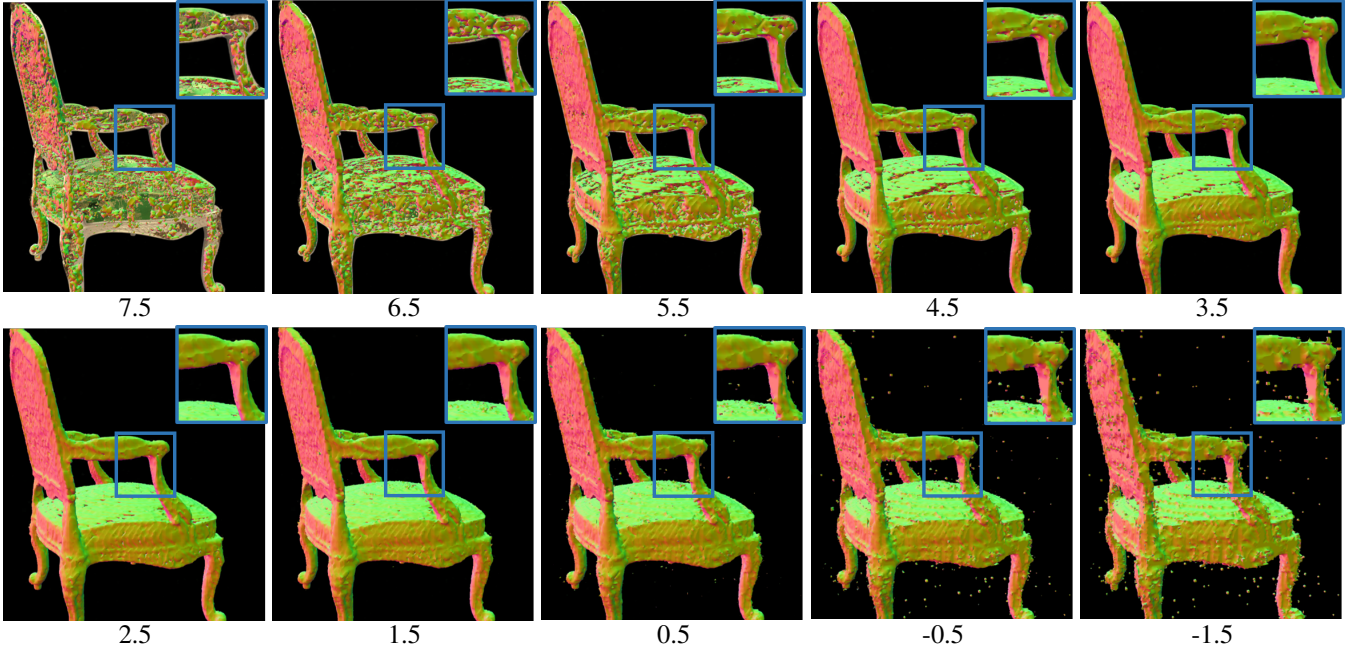


Figure 3. **Simultaneous visualization of extracted geometry normal map and corresponding rendering frame.**

NeRF-Synthetic Chair scene [3] from scratch, extract the two geometry proxies using marching cubes with thresholds -1.5 and 5.5 (refer to Figure 3 for geometry visualization). To ensure the fairness of experiments, we leverage the same sampled camera poses for distillation view synthesis as those used in the main paper, and the same configurations for optimizing neural duplex radiance fields. We present both quantitative results and qualitative results in Table 3 and Figure 2, which effectively demonstrate the generalizations of the neural duplex radiance fields.

Another Perspective to Understand Duplex Radiance Fields. To better understand the mechanism of the proposed neural duplex radiance field and how it works, we simultaneously visualize the geometry and the rendering frame under different thresholds from a pretrained Instant-NGP model [4] in Figure 3. We have several important observations: 1) Larger thresholds will lead to an under-estimated mesh, which cannot cover the ground-truth surface, but include abundant geometric details. 2)

Table 4. **Comparisons with different mesh layer settings.** Time: Rasterization overheads of the CUDA version.

Threshold	A ¹	B ¹	C ¹	Ours ²	D ³	E ⁴	F ⁴
5×10^{-5}							✓
1×10^{-4}		✓		✓	✓	✓	✓
5×10^{-4}						✓	
1×10^{-3}	✓				✓		
5×10^{-3}						✓	
1×10^{-2}			✓	✓	✓	✓	✓
5×10^{-2}							✓
PSNR ↑	27.52	25.50	27.05	32.67	32.84	33.03	32.90
SSIM ↑	0.941	0.921	0.935	0.975	0.976	0.977	0.976
LPIPS ↓	0.068	0.088	0.069	0.024	0.022	0.021	0.022
Time (ms) ↓	2.41	2.59	2.03	4.64	7.23	9.64	8.85

Smaller thresholds will lead to an over-estimated mesh with noise, but it can effectively wrap up the underlying true surface and contains fewer holes. 3) The distribution of the geometries is *spatially heterogeneous*. For example, the triangle surfaces extracted from threshold in the range [2.5, 6.5] are not concurrently inside or outside of the scenes. In one mesh, some regions may be covered and others may not. Hence, based on these observations, it will generally be unlikely that a single surface extracted from a NeRF can represent the 3D scene without sacrificing rendering quality. In contrast, our proposed neural duplex radiance field can effectively resolve these issues through learning the combination of two-layer duplex meshes while achieving high-quality rendering.

Will more mesh layers lead to better results? As shown in Table 4, increasing the mesh layers can further increase rendering quality, but the gain is **minor** compared to the improvement from a single mesh (‘A’/‘B’/‘C’) to our two-layer duplex mesh (‘Ours’). Meanwhile, increasing the mesh layers will result in linear growth of memory footprint and rasterization time, since the hardware rasterization cannot run in parallel for multiple meshes, which will negatively impact the real-time application. Hence, our duplex radiance field provides the optimal trade-off between quality and speed.

References

- [1] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. TensorRF: Tensorial radiance fields. In *ECCV*, 2022.
- [2] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *CVPR*, pages 1251–1258, 2017.
- [3] 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.
- [4] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–15, 2022.