

Neural Lumigraph Rendering

Petr Kellnhofer^{1,2}, Lars C. Jebe¹, Andrew Jones¹, Ryan Spicer¹, Kari Pulli¹, Gordon Wetzstein^{2,1}

¹Raxium ²Stanford University

{pkellnhofer, ljebe, ajones, rspicer, kpulli, gwetzstein}@raxium.com

Abstract

Novel view synthesis is a challenging and ill-posed inverse rendering problem. Neural rendering techniques have recently achieved photorealistic image quality for this task. State-of-the-art (SOTA) neural volume rendering approaches, however, are slow to train and require minutes of inference (i.e., rendering) time for high image resolutions. We adopt high-capacity neural scene representations with periodic activations for jointly optimizing an implicit surface and a radiance field of a scene supervised exclusively with posed 2D images. Our neural rendering pipeline accelerates SOTA neural volume rendering by about two orders of magnitude and our implicit surface representation is unique in allowing us to export a mesh with view-dependent texture information. Thus, like other implicit surface representations, ours is compatible with traditional graphics pipelines, enabling real-time rendering rates, while achieving unprecedented image quality compared to other surface methods. We assess the quality of our approach using existing datasets as well as high-quality 3D face data captured with a custom multi-camera rig.

1. Introduction

Novel view synthesis and 3D shape estimation from 2D images are inverse problems of fundamental importance in applications as diverse as photogrammetry, remote sensing, visualization, AR/VR, teleconferencing, visual effects, and games. While traditional 3D computer vision pipelines have been studied for decades, only emerging neural rendering techniques have been able to achieve photorealistic quality for novel view synthesis (e.g., [38, 56]).

State-of-the-art neural rendering approaches, such as neural radiance fields [38], however, do not offer real-time framerates, which severely limits their applicability to the aforementioned problems. This limitation is primarily imposed by the choice of implicit neural scene representation and rendering algorithm, namely a volumetric representation that requires a custom neural volume renderer. Neural

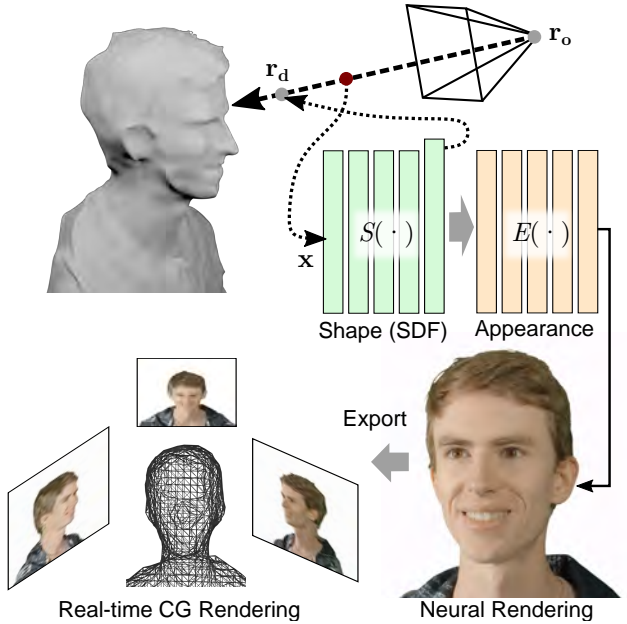


Figure 1. Overview of our framework. Given a set of multi-view images, we optimize representation networks modeling shape and appearance of a scene end to end using a differentiable sphere tracer. The resulting models can be exported to enable view-dependent real-time rendering using traditional graphics pipelines.

surface representations, for example using signed distance functions (SDFs) [41, 18, 2, 60], occupancy fields [35], or feature-based representations [54], on the other hand implicitly model the surface of objects, which can be extracted using the marching cubes algorithm [32] and exported into traditional mesh-based representations for real-time rendering. Although implicit neural surface representations have recently demonstrated impressive performance on shape reconstruction [60], their performance on view interpolation and synthesis tasks is limited. Thus, SOTA neural rendering approaches either perform well for view synthesis [38] or 3D shape estimation [60], but not both.

Here, we adopt an SDF-based sinusoidal representation

network (SIREN) as the backbone of our neural rendering system. While these representations have recently demonstrated impressive performance on representing shapes via direct 3D supervision with point clouds [52], we are the first to demonstrate how to leverage SIREN’s extreme capacity in the context of learning 3D shapes using 2D supervision with images via neural rendering. For this purpose, we devise a novel loss function that maintains SIREN’s high-capacity encoding for the supervised images while constraining it in the angular domain to prevent overfitting on these views. This training procedure allows us to robustly fit a SIREN-based SDF directly to a sparse set of multi-view images. Our 2D-supervised implicit neural scene representation and rendering approach performs on par with NeRF on view interpolation tasks while providing a high-quality 3D surface that can be directly exported for real-time rendering at test time.

Specifically, we make the following contributions:

- We develop a neural rendering framework comprising an implicit neural 3D scene representation, a neural renderer, and a custom loss function for training. This approach achieves 10× higher rendering rates than NeRF while providing comparable, SOTA image quality with the additional benefit of optimizing an implicitly defined surface.
- We demonstrate how both shape and view-dependent appearance of our neural scene representation can be exported and rendered in real time using traditional graphics pipelines.
- We also build a custom camera array and capture several datasets of faces and heads for evaluating our approach and baselines. These data are available on the project website.¹

2. Related Work

Traditional 3D computer vision pipelines use structure-from-motion and multi-view-stereo algorithms to estimate sparse point clouds, camera poses, and textured meshes from 2D input views (e.g., [55, 49, 7, 47, 48]). Re-rendering these scene representations, however, does not achieve photorealistic image quality. As an alternative, image-based rendering techniques have been explored for decades [50]. Lumigraph rendering [17, 4] stands out among these methods as an approach that leverages proxy scene geometry to interpolate the captured views better. Still, these traditional approaches have not demonstrated photorealistic view synthesis for general 3D scenes.

Emerging neural scene representations often model an object or scene explicitly using some 3D proxy geometry,

such as an imperfect mesh [21, 57, 45, 62] or depth map [44] estimated by multi-view stereo or other means, an object-specific shape template [25], a multi-plane [64, 37, 11] or multi-sphere [3, 1] image, or a volume [53, 31]. An overview of recent neural rendering techniques, including extensive discussions of explicit representations, is provided by Tewari et al. [56].

As opposed to explicit representations, emerging neural implicit scene representations promise 3D-structure-aware, continuous, memory-efficient representations for shape parts [15, 14], objects [41, 36, 2, 18, 60, 8, 5], or scenes [10, 54, 23, 43, 52]. These representations implicitly define an object or a scene using a neural network and can be supervised directly with 3D data, such as point clouds, or with 2D multi-view images [46, 54, 40, 39, 38, 60, 29, 24, 30, 27]. It is important to distinguish between neural implicit representations that use implicitly defined volumes [38, 29, 28] from those using implicitly defined surfaces, for example represented as signed distance functions (SDFs) [41, 36, 2, 18, 54, 23, 43] or occupancy networks [35, 6, 39]. Surface-based representations allow for traditional mesh representations to be extracted and rendered efficiently with traditional computer graphics pipelines.

The neural rendering methods closest to our work are neural radiance fields (NeRF) [38], which provide the best image quality for view synthesis to date but do not directly model object shape, and implicit differentiable renderer (IDR) [60], which recently demonstrated SOTA performance for shape estimation but which does not achieve the same quality as NeRF for view synthesis. Our work leverages emerging sinusoidal representation networks (SIREN) [52] to achieve both of these capabilities simultaneously. This is crucial for the neural rendering pipeline we propose, which first learns a neural implicit surface representation and then exports it into a format that is compatible with existing real-time graphics pipelines. Although SIRENs have been proposed in prior work [52], we are the first to demonstrate how to leverage their impressive capacity with 2D multi-view image supervision – a non-trivial task due to SIREN’s extreme overfitting behavior.

3. Neural Rendering Pipeline

In this section, we describe our differentiable neural rendering pipeline, which is illustrated in Fig. 1.

3.1. Representation

We represent both shape and appearance of 3D objects using implicit functions in a framework similar to IDR [60]. Unlike previous work, however, our network architecture builds on sinusoidal representation networks (SIREN) [52], which allow us to represent signals of significantly higher complexity within the same number of learnable param-

¹<http://www.computationalimaging.org/publications/nlr/>

ters compared with common non-periodic multilayer perceptrons (MLP).

We express the continuous shapes of a scene as the zero-level set $S_0 = \{\mathbf{x} \mid S(\mathbf{x}) = 0\}$ of a signed distance function (SDF)

$$S(\mathbf{x}; \theta) : \mathbb{R}^3 \rightarrow \mathbb{R}, \quad (1)$$

where $\mathbf{x} \in \mathbb{R}^3$ is a location in 3D space and θ are the learnable parameters of our SIREN-based SDF representation.

Next, we model appearance as a spatially varying emission function, or radiance field, E for directions $\mathbf{r}_d \in \mathbb{R}^3$ defined in a global coordinate system. This formulation does not allow for relighting but it enables photorealistic reconstruction of the appearance of a scene under fixed lighting conditions. We leave the problem of modeling lighting and shading as an avenue of future work.

We additionally condition E by the local normal direction $\mathbf{n} = \nabla_{\mathbf{x}} S(\mathbf{x})$ as computed by automatic differentiation. This does not constrain any degrees of freedom but it has been shown to improve the training performance [60]. Finally, we also reuse θ to increase the network capacity and allow for modeling of fine spatial details and micro-reflections that are of a notably higher spatial complexity than the underlying shape. Together, we express the radiance field as

$$E(\mathbf{x}, \mathbf{r}_d, \mathbf{n}; \theta, \phi) : \mathbb{R}^9 \rightarrow \mathbb{R}^3, \quad (2)$$

to represent RGB appearance using the additional learnable parameters ϕ .

3.2. Neural Rendering

The goal of neural rendering is to project a 3D neural scene representation into one or multiple 2D images. We solve this task in two steps: 1) We find the 3D surface as the zero-level set S_0 closest to the camera origin along each ray; 2) We resolve the appearance by sampling the local radiance E .

To address 1), we sphere trace the SDF to find S_0 [20]. For this purpose, we define a view and a projection matrix, $V \in \mathbb{R}^{4 \times 4}$ and $P \in \mathbb{R}^{4 \times 4}$, similar to OpenGL’s rendering API [59]. A ray origin \mathbf{r}_o and direction \mathbf{r}_d for an output pixel at relative projection plane location $\mathbf{u} \in [-1, 1]^2$ is then

$$\mathbf{r}_o = (V^{-1} \cdot [0, 0, 1, 0]^T)_{x,y,z}, \quad (3)$$

$$\mathbf{r}_d = \nu((P \cdot V)^{-1} \cdot [\mathbf{u}_x, \mathbf{u}_y, 0, 1]^T), \quad (4)$$

where $(\cdot)_{x,y,z}$ are vector components and $\nu(\omega) = \omega_{x,y,z} / \|\omega_{x,y,z}\|$ is vector normalization.

The sphere-tracing algorithm minimizes $|S(\mathbf{x}, \theta)|$ along each ray using iterative updates of the form

$$\mathbf{x}_0 = \mathbf{r}_o, \quad \mathbf{x}_{i+1} = \mathbf{x}_i + S(\mathbf{x}_i) \mathbf{r}_d. \quad (5)$$

Finally, $S_0 = \{\mathbf{x}_n \mid S(\mathbf{x}_n) = 0\}$ is the zero-set of rays converged to a foreground object for the step count $n = 16$. A small residual $|S(\mathbf{x}_n)| < 0.005$ is tolerated in practice. As proposed in recent work [60, 24, 30], we only retain gradients in the last step rather than for all steps of the sphere tracer. This approach makes sphere tracing memory efficient. Please refer to the supplemental materials for additional details.

The appearance is directly sampled from our radiance field as $E(S_0, \mathbf{r}_d, \nabla S(S_0); \theta, \phi)$.

3.3. Loss Function

We supervise our 3D representation using a set of m multi-view 2D images $\mathbf{I} = \mathbb{R}^{m \times w \times h \times 3}$ with known object masks $\mathbf{M} = \mathbb{R}^{m \times w \times h}$ where 1 marks foreground. Our unique approach to leveraging SIREN as a neural representation in this setting is challenging, because of SIREN’s tendency to overfit the signal to the supervised views.

In total, we use four different constraints to optimize the end-to-end representation using mini-batches of image pixels \mathbf{U} with RGB values $\mathbf{I}_{\mathbf{U}}$ and object masks $\mathbf{M}_{\mathbf{U}}$.

First, we minimize an $L1$ image reconstruction error for the true foreground pixels $\mathbf{U}_f = \mathbf{U} \cap S_0 \cap \{\mathbf{U} \mid \mathbf{M}_{\mathbf{U}} = 1\}$ as

$$\mathcal{L}_R = \frac{1}{|\mathbf{U}|} \sum_{\mathbf{c} \in \mathbf{I}_{\mathbf{U}_f}} |E(\mathbf{x}, \mathbf{r}_d, \mathbf{n}; \theta, \phi) - \mathbf{c}|, \quad (6)$$

where \mathbf{c} is an RGB value of a foreground pixel in a mini-batch. Both $L1$ and $L2$ work well but we have found $L1$ to produce marginally sharper images.

Second, we regularize the S by an eikonal constraint

$$\mathcal{L}_E = \frac{1}{|\mathbf{U}|} \sum_{\mathbf{x}_r} \|(\|\nabla_{\mathbf{x}} S(\mathbf{x}_r; \theta)\|_2 - 1)\|_2^2 \quad (7)$$

to enforce its metric properties important for efficient sphere tracing [19, 24, 52, 60]. Random points \mathbf{x}_r are uniformly sampled from a cube which encapsulates the object’s bounding unit radius sphere.

Third, we restrict the coarse shape by enforcing its projected pattern to fall within the boundaries of the object masks. For this purpose, we adopt the soft mask loss proposed in [60] defined for the non-foreground pixels and softness parameter α as

$$\mathcal{L}_M = \frac{1}{\alpha |\mathbf{U}|} \sum_{m \in \mathbf{M}_{\mathbf{U} \setminus \mathbf{U}_f}} \text{BCE}(\text{sigmoid}(-\alpha S_{min}), m), \quad (8)$$

where BCE is the binary cross entropy and $S_{min} = \arg \min_t S(\mathbf{r}_o + t \mathbf{r}_d; \theta)$ is the minimum S value along the entire ray approximated by dense sampling of t .

Finally, we regularize the radiance field E to avoid overfitting to training views. SIRENs have a remarkable regressive potential, which biases them to overfit the appearance

to the training views. We leverage this power to allow for encoding of photorealistic surface details, but we need to restrict the behavior of the E in the angular domain conditioned by \mathbf{r}_d to achieve favorable interpolation behavior. Inspired by multi-view projective texture mapping [9], we linearize the angular behavior using a smoothness term

$$\mathcal{L}_S = \frac{1}{|\mathbf{U}|} \sum \|\nabla_{\mathbf{r}_d}^2 E(\mathbf{x}, \mathbf{r}_d, \mathbf{n}; \theta, \phi)\|_2^2. \quad (9)$$

Note that such level of control is unique to SIREN and related architectures as they are C^∞ differentiable.

Together, we optimize parameters θ and ϕ as

$$\arg \min_{\theta, \phi} \mathcal{L}_R + w_E \mathcal{L}_E + w_M \mathcal{L}_M + w_S \mathcal{L}_S, \quad (10)$$

with weights $w_E = 0.1$, $w_M = 100$, and $w_S = 0.01$ for all of our experiments. We have not found the performance to be very sensitive to this choice with the exception of w_S where large values cause high-frequency artifacts in S .

3.4. Additional Training Details

We optimize the loss in mini-batches of 50,000 individual rays sampled uniformly across the entire training dataset. We have found a large batch size and uniform ray distribution to be critical to prevent local overfitting of SIREN, especially for the high-frequency function E .

We implement the MLPs representing S and E as SIRENs with 5 layers using 256 hidden units each. Additionally, we use Fourier features $\{\sin(2k\pi\mathbf{r}_d), \cos(2k\pi\mathbf{r}_d) \mid k \in 1 \dots 4\}$ in E to further support angular resolution [38, 60]. This strategy is necessary to fit the sparsely supervised rays well while \mathcal{L}_S enhances interpolation between them.

We initialize S to a unit sphere of radius 0.5 by pre-training to a procedural shape as described in [52]. We trace the object rays in a larger sphere of radius 1, but we have found that the smaller initial radius improves the initial fit as well as the consequent convergence rate.

We implement our method in PyTorch [42] and optimize the loss using the Adam solver [26] with an initial learning rate of 10^{-4} decreased by a factor of 2 every 40,000 batches for the overall training length of 150,000 batches on a single Nvidia GPU RTX 2080Ti.

4. Real-time Rendering Pipeline

While we show that SIREN is remarkably efficient in shape and appearance representation with 2D supervision, the required sphere tracer does not run at real-time rates for moderate to high image resolutions. To overcome this challenge, we show how to leverage the compactness of our surface-based representation and convert our neural model to a triangular mesh suitable for real-time computer graphics applications. For this purpose, we leverage unstructured

lumigraph rendering, which preserves view-dependent effects learned by our neural representation [4].

4.1. Mesh extraction

First, we use the marching cubes algorithm [33] to extract a high-resolution surface mesh from the SDF S voxelized at a resolution of 512^3 . Instead of extracting the zero-level set, we found that offsetting the iso-surface of S by 0.5% of the object radius in the outside direction optimizes the resulting image quality in practice. To export the appearance, we resample the optimized emissivity function E to synthesize projective textures \mathbf{T}_i for N camera poses and corresponding projection matrices. The ability to resample the camera poses for efficient viewing space coverage is a key feature of our method and we explore the choice of N and camera distributions in the supplement.

4.2. Rendering

First, we rasterize the extracted mesh using OpenGL [59] and project the vertex positions to each pixel. Next, we compute angles $\tau_{1 \dots N}$ between the ray towards the current rendering camera and the rays towards each of the N projective texture map viewpoints. We then apply the unstructured lumigraph rendering technique of Buehler et al. [4] to blend contributions from the first $k = 5$ textures, sorted by τ_i in ascending order, yielding a rendered image

$$\mathbf{R} = \sum_{i=1 \dots k} w_i \mathbf{T}_i, \quad (11)$$

where the weights w_i are computed as

$$\hat{w}_i = 1/\tau_i(1 - \tau_i/\tau_k), \quad (12)$$

$$w_i = \hat{w}_i / \sum_{i=1 \dots k} \hat{w}_i. \quad (13)$$

This formulation satisfies the epipolar consistency by converging to an exclusive mapping by texture \mathbf{T}_j when $\tau_j \rightarrow 0$ [4]. Additionally, we discard samples from occluded textures by setting their w_i to zero. Occlusions are detected by a comparison between the pre-rendered depth associated with a texture and the distance between the mesh voxel and the texture viewpoint. The same technique is commonly used in real-time graphics for shadow mapping.

4.3. Evaluation

Efficiency We compare the efficiency of our real-time rasterized neural lumigraph renderer (NLR-RAS) with our sphere-traced renderer (NLR-ST, Sec. 3) along with other baselines in Table 1. We observe, that although both NLR-ST and IDR are based on sphere tracing, the capacity of SIREN allows for a smaller and faster model, which is evident by the model size. Furthermore, the results show how costly the implicit volumetric rendering is. In conclusion, only the explicit representations of Colmap and our NLR-RAS allow for truly real-time performance with framerates over 60 fps at HD resolution on commodity hardware.

Method	Render time [s]	Model size [MB]
Colmap [48]	Real-time	30.39
IDR [60]	45	11.13
NV [31]	0.65	438.36
NeRF [37]	150	2.27
NLR-ST	13	2.07
NLR-RAS	Real-time	34.68

Table 1. Rendering time and representation size comparison for the DTU scan 65 [22] at 1600×1200 pixel resolution. “Real-time” denotes framerates of at least 60 fps.



Figure 2. Our custom camera array comprising 16 GoPro HERO7 and 6 central Back-Bone H7PRO cameras (large circular lenses).

Image quality Both the quantitative comparisons in Tables 2, 3 and qualitative examples in Figures 3, 4 demonstrate the high NLR-RAS rendering quality. While lower than that of our NLR-ST renderer, the NLR-RAS still achieves PSNRs far superior to other explicit (Colmap) and implicit (IDR) surface representations.

5. Camera Array and Data

Human Head Video Dataset Our dataset consists of 7 multiview captures showing a person performing facial expressions.

Camera Array Our custom camera array that was used to capture the dataset is comprised of 16 GoPro HERO7 action cameras and 6 Back-Bone H7PRO cameras. The Back-Bone cameras are modified GoPro cameras that can fit a standard C-Mount lens. Compared to the unmodified GoPro cameras, the Back-Bone cameras have a narrower field-of-view (FoV) and are thus able to capture the subject in more detail. We capture at 4k / 30 fps in portrait orientation with the Back-Bone cameras and at 1080p / 60 fps in landscape orientation with the GoPro cameras. Figure 2 shows a frontal view of the camera array with the six Back-Bone cameras in the center of the array and the GoPro cameras placed around them. We capture our subjects from 60 cm distance and cover approximately 100° .

Synchronization We trigger the camera shutter with a WiFi remote. The cameras do not support a generator lock, so during capture they are only loosely synchronized. We always capture videos for our dataset, even in the cases in which we only use a static frame. To improve synchronization, we flash an ArUco marker [13] on a cellphone before

each take. We then detect the first frame that sees the marker in each video which allows us to synchronize the cameras with an accuracy of 1 frame or better.

6. Experiments

In this section we show that our method is able to achieve state-of-the-art image reconstruction quality on-par with volumetric methods such as NeRF [38] while allowing for efficient surface reconstruction utilized for real-time rendering in Sec. 4.

Baselines We compare our method to novel view synthesis techniques with various scene representations. Specifically, we compare to the traditional multi-view stereo of Colmap [48], the explicit volumetric representation of Neural Volumes (NV) [31], the implicit volume representation of NeRF [38], and the implicit signed distance function of IDR [60]. Please refer to the supplemental material for implementation details.

Regression performance We have used the popular DTU MVS dataset [22] with 49 or 64 calibrated camera images along with object masks provided by previous work [60, 39] to measure the image reconstruction error metrics. We held out three views for testing. Table 2 shows that our method achieves SOTA training error comparable with NeRF. Moreover, our image quality is significantly better than that of our closest competitor, IDR. We attribute this major separation to the unparalleled representation capacity of SIRENS. A qualitative comparison is available in Fig. 3.

Additionally, we report the shape reconstruction error as Chamfer distance from the ground-truth provided in the dataset. Although the shape reconstruction is not our explicit goal, we note that our error is on par with other techniques, though worse than IDR which explicitly focuses on this problem [60]. We observe that this emerges as a trade-off between the accuracy of view-dependent and high-frequency details in the image reconstruction on one hand, and the view consistency reflected in the geometry on the other one.

View interpolation Our angular smoothness loss \mathcal{L}_S is specifically designed to avoid collapse of the emissivity function E for interpolated views. We tested its efficiency quantitatively by measuring the image reconstruction error on test views the held out from results in Table 2. Please refer to the supplement for details. As expected, there is a measurable quality drop when compared to the training views observed consistently for all of the methods. However, the interpolated views produced by our method maintain many of the favorable characteristics from the regres-

Scan	Colmap [48]				IDR [60]				NeRF [37]				NLR-ST				NLR-RAS		
	CD↓	PSNR↑	SSIM↑	LPIPS↓	CD↓	PSNR↑	SSIM↑	LPIPS↓	CD↓	PSNR↑	SSIM↑	LPIPS↓	CD↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓
65	2.42	15.25	0.862	0.143	0.70	23.87	0.948	0.094	<u>1.02</u>	33.57	0.962	0.055	<u>1.03</u>	<u>32.13</u>	<u>0.961</u>	<u>0.063</u>	31.46	0.960	0.066
97	0.72	11.93	0.843	0.154	<u>1.09</u>	23.02	0.921	0.117	1.43	28.28	0.933	<u>0.089</u>	1.16	28.48	0.939	0.088	<u>28.37</u>	<u>0.939</u>	<u>0.089</u>
106	<u>0.78</u>	15.75	0.887	0.141	0.58	20.97	0.907	0.183	0.84	33.50	0.947	0.092	<u>0.82</u>	<u>32.98</u>	0.951	0.083	31.32	<u>0.951</u>	<u>0.090</u>
118	0.44	22.73	0.921	0.091	<u>0.50</u>	22.62	0.944	0.139	0.88	35.62	0.966	0.070	<u>1.71</u>	<u>34.87</u>	<u>0.964</u>	<u>0.075</u>	33.33	0.963	0.078

Table 2. Image error metrics PSNR, SSIM [58] and LPIPS [61] computed on DTU [22] for the supervised views. See Supp. for metrics of the held-out views. The Chamfer distance (CD) is computed based on the scripts from [22]. Best scores in bold, second best underlined.



Figure 3. Reconstructed shape and image from various multi-view datasets. The reprojection error is listed in dB of PSNR.

sion case. We provide a qualitative comparison of both supervised and interpolated views in Fig. 4 and in our video.

Human representation View-synthesis of human subjects is particularly challenging due to the complex reflection properties of skin, eyes and hair, as well as a lack of high-quality multi-view data. We address the first challenge with our high-capacity representation network and the lat-

ter with our own dataset described in Sec. 5. Additionally, we provide experimental results for 360 degree human captures provided by Volucap GmbH [16] and high-resolution face captures from the Digital Ira project [12]. Refer to the supplemental for a detailed description of these data. Table 3 summarizes the reconstruction errors and Fig. 3 shows a few example scenes. Similar trends as in the DTU datasets can again be observed. Interestingly, our method achieves

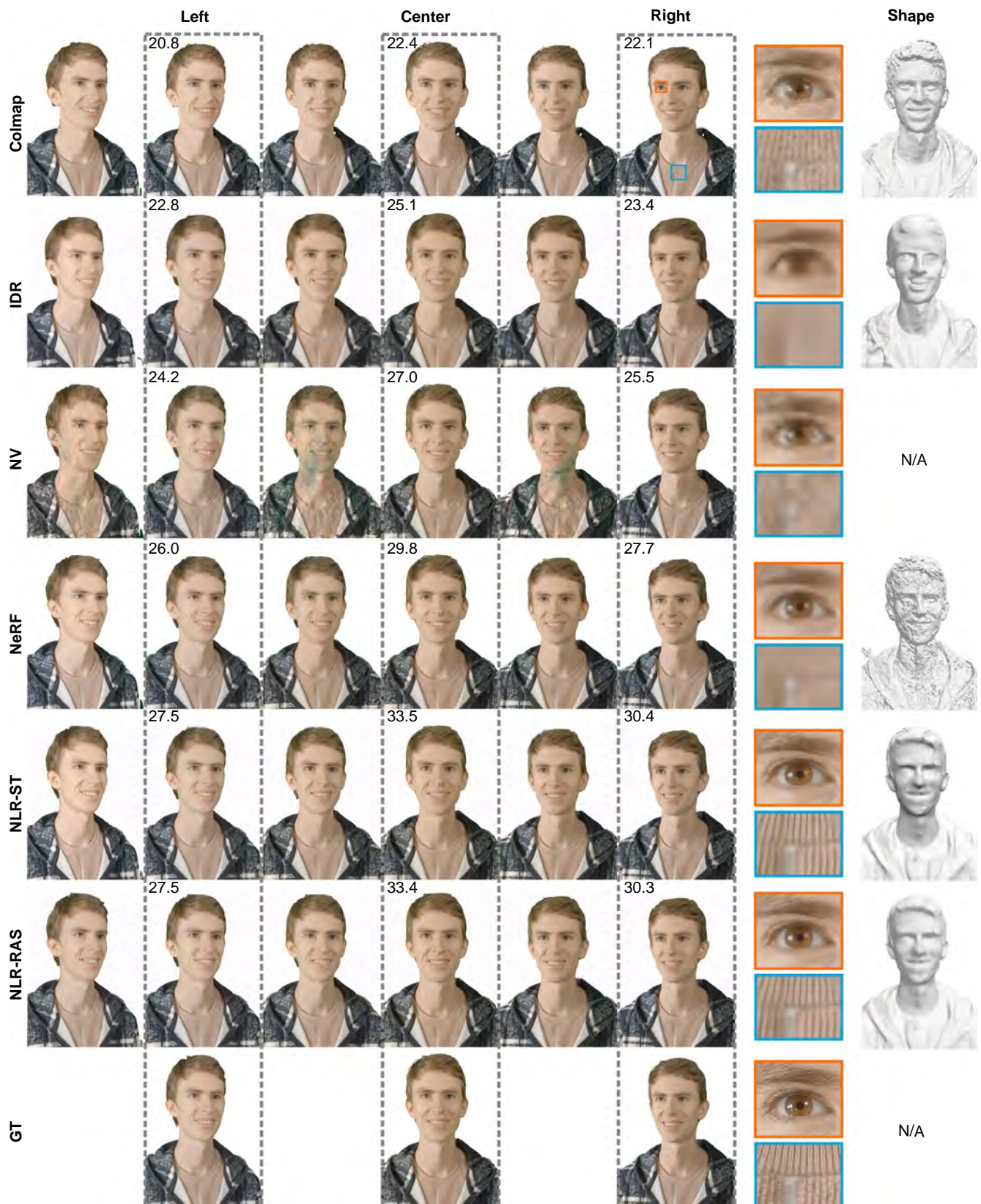


Figure 4. A qualitative comparison for a frame in our dataset. Three supervised views (dashed outline) along with three interpolated views and learned shape. Close-ups enhance detail of the rightmost view. PSNR showed for supervised views.



Figure 5. Learned shapes with interpolated close-ups (left) and reconstruction metrics for each condition (right) in our ablation study.

Dataset	Colmap	IDR	NV	NeRF	NLR-ST	N.-RAS
Volucap (1)	19.6/.037	22.3/.043	29.3/.034	32.7/.026	28.6/.022	28.6/.022
Dig. Ira (1)	Fail.	23.9/.286	26.5/.287	31.2/.267	31.6/.255	31.1/.260
Ours (7)	17.6/.187	22.3/.202	25.1/.186	28.5/.171	30.5/.147	30.3/.151

Table 3. Average reconstruction PSNR/LPIPS [61] scores computed across datasets (number of scenes in parentheses). *Italic*: Only 5 scenes tested. See Supplement for an extended version.

a bigger advantage for very high-resolution (3000×4000 px) detailed images in our own dataset. We speculate that this shows that the traditional ReLU based networks used by IDR and NeRF have reached their capacity, while the explicit representations of Colmap and NV lack easy scaling. Once again, this does not come at cost of interpolation properties as shown in Figure 4 and our videos.

Ablation study Finally, we verify that the performance of our method is based on the choice of our representation and training procedure. In Figure and Table 5, we compare several variants of our method on the scene in Fig. 4. A standard MLP with ReLU does not have the capacity to train a detailed representation (1). SIREN remedies this but tends to quickly overfit to the trained pixels (2). We resolve this first by adding our angular smoothness loss \mathcal{L}_S that regularizes behavior in the angular domain (3), and then by increasing the batch size in order to achieve spatially uniform image quality (4). Additional Fourier Features [38] for the ray direction remove low frequency noise in E (5).

7. Discussion

In summary, we propose a neural rendering framework that optimizes an SDF-based implicit neural scene representation given a set of multi-view images. This framework is unique in combining a representation network architecture using periodic activations with a sphere-tracing-based neural renderer that estimates the shape and view-dependent appearance of the scene. Enabled by a novel loss function that is applied during training, our framework achieves a very high image quality that is comparable with state-of-the-art novel view synthesis methods. As opposed to those methods, our neural representation can be directly

	Act.	\mathcal{L}_S	Batch size	Extra FF	PSNR \uparrow		LPIPS \downarrow [61]	
					Train	Test	Train	Test
(1)	Relu	No	2K	No	24.99	10.57	0.171	0.285
(2)	Sine	No	2K	No	28.49	19.27	0.161	0.239
(3)	Sine	Yes	2K	No	28.26	24.77	0.151	0.179
(4)	Sine	Yes	50K	No	30.06	25.20	0.143	0.172
(5)	Sine	Yes	50K	Yes	30.75	26.04	0.132	0.151

converted into a mesh with view-dependent textures that enable high-quality 3D image synthesis in real time using traditional graphics pipelines.

Our approach is not without limitations. Currently, we only consider emissive radiance functions that are adequate to model a scene under fixed lighting conditions. Future work could additionally consider dynamic lighting and shading, which some recent neural rendering approaches have started to incorporate [34, 62]. Further, similar to IDR [60], our method requires annotated object masks. Automatic image segmentation could be explored in the future to address this. Although the synthesized image quality of our approach is competitive with the state of the art, the proxy shapes produced by our method are not quite as accurate as alternative approaches [48, 60]. While this is not important for the novel view synthesis application we consider in this paper, other applications may benefit from estimating more accurate shapes. This includes occasionally visible seam artifacts caused by inaccuracies of the camera calibration. Similar to some other recent neural rendering pipelines, ours focuses on overfitting a neural representation on a single 3D scene. An interesting avenue of future work includes the learning of shape spaces, or priors, for certain types of objects, such as faces. While several methods have explored related strategies using conditioning-by-concatenation [41, 35], hypernetwork [54], or meta-learning [51] approaches using synthetic data, there is a lack of publicly available photorealistic multi-view image data. In hope of mitigating this shortcoming, we released our datasets on the project website. Still, these data may be insufficiently large for learning priors. Finally, although the inference time of our method is fast, the training time is still slow. This hurdle along with the limited computational resources at our disposal is the primary reason preventing us from exploring dynamic video sequences.

Conclusion Emerging neural rendering approaches are starting to outperform traditional vision and graphics approaches. Yet, traditional graphics pipelines still offer significant practical benefits, such as real-time rendering rates, over these neural approaches. With our work, we take a significant step towards closing this gap, which we believe to be a critical aspect for making neural rendering practical.

Acknowledgments

We would like to thank Volucap GmbH for providing a multi-view video captured using their studio setup as well USC Institute for Creative Technologies for providing a sample of the Digital Ira dataset.

References

- [1] Benjamin Attal, Selena Ling, Aaron Gokaslan, Christian Richardt, and James Tompkin. MatryODShka: Real-time 6DoF video view synthesis using multi-sphere images. In *Proc. ECCV*, Aug. 2020. 2
- [2] Matan Atzmon and Yaron Lipman. Sal: Sign agnostic learning of shapes from raw data. In *Proc. CVPR*, 2020. 1, 2
- [3] Michael Broxton, John Flynn, Ryan Overbeck, Daniel Erickson, Peter Hedman, Matthew Duvall, Jason Dourgarian, Jay Busch, Matt Whalen, and Paul Debevec. Immersive light field video with a layered mesh representation. *ACM Trans. Graph. (SIGGRAPH)*, 39(4), 2020. 2
- [4] Chris Buehler, Michael Bosse, Leonard McMillan, Steven Gortler, and Michael Cohen. Unstructured lumigraph rendering. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 425–432, 2001. 2, 4
- [5] Rohan Chabra, Jan Eric Lenssen, Eddy Ilg, Tanner Schmidt, Julian Straub, Steven Lovegrove, and Richard Newcombe. Deep local shapes: Learning local sdf priors for detailed 3d reconstruction. *arXiv preprint arXiv:2003.10983*, 2020. 2
- [6] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proc. CVPR*, pages 5939–5948, 2019. 2
- [7] Alvaro Collet, Ming Chuang, Pat Sweeney, Don Gillett, Dennis Evseev, David Calabrese, Hugues Hoppe, Adam Kirk, and Steve Sullivan. High-quality streamable free-viewpoint video. *ACM Trans. Graph. (SIGGRAPH)*, 34(4), 2015. 2
- [8] Thomas Davies, Derek Nowrouzezahrai, and Alec Jacobson. Overfit neural networks as a compact shape representation, 2020. 2
- [9] Paul Debevec, Yizhou Yu, and George Borshukov. Efficient view-dependent image-based rendering with projective texture-mapping. In *Eurographics Workshop on Rendering Techniques*, pages 105–116. Springer, 1998. 4
- [10] SM Ali Eslami, Danilo Jimenez Rezende, Frederic Besse, Fabio Viola, Ari S Morcos, Marta Garnelo, Avraham Ruderman, Andrei A Rusu, Ivo Danihelka, Karol Gregor, et al. Neural scene representation and rendering. *Science*, 360(6394):1204–1210, 2018. 2
- [11] John Flynn, Michael Broxton, Paul Debevec, Matthew Duvall, Graham Fyffe, Ryan Overbeck, Noah Snively, and Richard Tucker. Deepview: View synthesis with learned gradient descent. In *Proc. CVPR*, pages 2367–2376, 2019. 2
- [12] Graham Fyffe, Andrew Jones, Oleg Alexander, Ryosuke Ichikari, and Paul Debevec. Driving high-resolution facial scans with video performance capture. *ACM Trans. Graph.*, 34(1), Dec. 2015. 6, 14
- [13] Sergio Garrido-Jurado, Rafael Muñoz-Salinas, Francisco José Madrid-Cuevas, and Manuel Jesús Marín-Jiménez. Automatic generation and detection of highly reliable fiducial markers under occlusion. *Pattern Recognition*, 47(6):2280–2292, 2014. 5
- [14] Kyle Genova, Forrester Cole, Avneesh Sud, Aaron Sarna, and Thomas Funkhouser. Local deep implicit functions for 3d shape. In *Proc. CVPR*, 2020. 2
- [15] Kyle Genova, Forrester Cole, Daniel Vlasic, Aaron Sarna, William T Freeman, and Thomas Funkhouser. Learning shape templates with structured implicit functions. In *Proc. ICCV*, pages 7154–7164, 2019. 2
- [16] Volucap GmbH. A multiview capture data sample, 2020. 6, 14
- [17] Steven J. Gortler, Radek Grzeszczuk, Richard Szeliski, and Michael F. Cohen. The lumigraph. In *ACM SIGGRAPH*, page 43–54, 1996. 2
- [18] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. Implicit geometric regularization for learning shapes. In *Proc. ICML*, 2020. 1, 2
- [19] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. Implicit geometric regularization for learning shapes. In *Proceedings of Machine Learning and Systems 2020*, pages 3569–3579. 2020. 3
- [20] John C Hart. Sphere tracing: A geometric method for the antialiased ray tracing of implicit surfaces. *The Visual Computer*, 12(10):527–545, 1996. 3
- [21] Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel Brostow. Deep blending for free-viewpoint image-based rendering. *ACM Transactions on Graphics (TOG)*, 37(6):1–15, 2018. 2
- [22] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engil Tola, and Henrik Aanæs. Large scale multi-view stereopsis evaluation. In *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pages 406–413. IEEE, 2014. 5, 6, 13, 14, 15
- [23] Chiyu Jiang, Avneesh Sud, Ameesh Makadia, Jingwei Huang, Matthias Nießner, and Thomas Funkhouser. Local implicit grid representations for 3d scenes. In *Proc. CVPR*, pages 6001–6010, 2020. 2
- [24] Yue Jiang, Dantong Ji, Zhizhong Han, and Matthias Zwicker. Sdfdiff: Differentiable rendering of signed distance fields for 3d shape optimization. In *Proc. CVPR*, 2020. 2, 3, 15
- [25] Angjoo Kanazawa, Shubham Tulsiani, Alexei A Efros, and Jitendra Malik. Learning category-specific mesh reconstruction from image collections. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 371–386, 2018. 2
- [26] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 4
- [27] Amit Kohli, Vincent Sitzmann, and Gordon Wetzstein. Semantic implicit neural scene representations with semi-supervised training. In *Proc. 3DV*, 2020. 2
- [28] David B Lindell, Julien NP Martel, and Gordon Wetzstein. Autoint: Automatic integration for fast neural volume rendering. In *Proc. CVPR*, 2021. 2

- [29] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *NeurIPS*, 2020. 2
- [30] Shaohui Liu, Yinda Zhang, Songyou Peng, Boxin Shi, Marc Pollefeys, and Zhaopeng Cui. Dist: Rendering deep implicit signed distance function with differentiable sphere tracing. In *Proc. CVPR*, 2020. 2, 3, 15
- [31] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM Trans. Graph. (SIGGRAPH)*, 38(4), 2019. 2, 5, 14, 15
- [32] William E. Lorensen and Harvey E. Cline. Marching cubes: A high resolution 3d surface construction algorithm. In *ACM SIGGRAPH*, page 163–169, 1987. 1
- [33] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. *ACM siggraph computer graphics*, 21(4):163–169, 1987. 4
- [34] Abhimitra Meka, Christian Häne, Rohit Pandey, Michael Zollhöfer, Sean Fanello, Graham Fyffe, Adarsh Kowdle, Xueming Yu, Jay Busch, Jason Dourgarian, Peter Denny, Sofien Bouaziz, Peter Lincoln, Matt Whalen, Geoff Harvey, Jonathan Taylor, Shahram Izadi, Andrea Tagliasacchi, Paul Debevec, Christian Theobalt, Julien Valentin, and Christoph Rhemann. Deep reflectance fields: High-quality facial reflectance field inference from color gradient illumination. *ACM Trans. Graph. (SIGGRAPH)*, 38(4), 2019. 8
- [35] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proc. CVPR*, 2019. 1, 2, 8
- [36] Mateusz Michalkiewicz, Jhony K Pontes, Dominic Jack, Mahsa Baktashmotlagh, and Anders Eriksson. Implicit surface representations as layers in neural networks. In *Proc. ICCV*, pages 4743–4752, 2019. 2
- [37] Ben Mildenhall, Pratul P. Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Trans. Graph. (SIGGRAPH)*, 38(4), 2019. 2, 5, 6, 14
- [38] 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 *Proc. ECCV*, 2020. 1, 2, 4, 5, 8, 14
- [39] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proc. CVPR*, 2020. 2, 5, 14
- [40] Michael Oechsle, Lars Mescheder, Michael Niemeyer, Thilo Strauss, and Andreas Geiger. Texture fields: Learning texture representations in function space. In *Proc. ICCV*, 2019. 2
- [41] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proc. CVPR*, 2019. 1, 2, 8
- [42] 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 H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. 4
- [43] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *Proc. ECCV*, 2020. 2
- [44] Eric Penner and Li Zhang. Soft 3d reconstruction for view synthesis. *ACM Trans. Graph.*, 36(6), Nov. 2017. 2
- [45] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *Proc. ECCV*, 2020. 2
- [46] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *Proc. ICCV*, pages 2304–2314, 2019. 2
- [47] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4104–4113, 2016. 2
- [48] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016. 2, 5, 6, 8, 12, 14
- [49] Steven M. Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. A comparison and evaluation of multi-view stereo reconstruction algorithms. In *Proc. CVPR*, page 519–528, 2006. 2
- [50] Harry Shum and Sing Bing Kang. Review of image-based rendering techniques. In *Visual Communications and Image Processing 2000*, volume 4067, pages 2–13. International Society for Optics and Photonics, 2000. 2
- [51] Vincent Sitzmann, Eric R. Chan, Richard Tucker, Noah Snavely, and Gordon Wetzstein. Metasdf: Meta-learning signed distance functions. In *NeurIPS*, 2020. 8
- [52] Vincent Sitzmann, Julien N.P. Martel, Alexander W. Bergman, David B. Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. In *Proc. NeurIPS*, 2020. 2, 3, 4
- [53] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhöfer. Deepvoxels: Learning persistent 3d feature embeddings. In *Proc. CVPR*, 2019. 2
- [54] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *Proc. NeurIPS 2019*, 2019. 1, 2, 8
- [55] Richard Szeliski. *Computer vision: algorithms and applications*. Springer, 2011 edition, 2010. 2
- [56] Ayush Tewari, Ohad Fried, Justus Thies, Vincent Sitzmann, Stephen Lombardi, Kalyan Sunkavalli, Ricardo Martin-Brualla, Tomas Simon, Jason Saragih, Matthias Nießner, et al. State of the art on neural rendering. *Proc. Eurographics*, 2020. 1, 2

- [57] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *ACM Transactions on Graphics (TOG)*, 38(4):1–12, 2019. 2
- [58] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 6, 14
- [59] Mason Woo, Jackie Neider, Tom Davis, and Dave Shreiner. *OpenGL programming guide: the official guide to learning OpenGL, version 1.2*. Addison-Wesley Longman Publishing Co., Inc., 1999. 3, 4
- [60] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Ronen Basri, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. In *Proc. NeurIPS*, 2020. 1, 2, 3, 4, 5, 6, 8, 12, 14, 15
- [61] 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. 6, 8, 14
- [62] Xiuming Zhang, Sean Fanello, Yun-Ta Tsai, Tiancheng Sun, Tianfan Xue, Rohit Pandey, Sergio Orts-Escolano, Philip Davidson, Christoph Rhemann, Paul Debevec, Jonathan T. Barron, Ravi Ramamoorthi, and William T. Freeman. Neural light transport for relighting and view synthesis, 2020. 2, 8
- [63] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3D: A modern library for 3D data processing. *arXiv:1801.09847*, 2018. 15
- [64] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. *ACM Trans. Graph. (SIGGRAPH)*, 37(4), 2018. 2

A. Real-time rendering analysis

In Section 4.3 of the paper, we compared performance of our real-time renderer to the neural renderer. Here, we complement this comparison by demonstrating that the rendering quality is a joint product of both the shape representation in S and the emissivity function E . To that goal, we use our new head image dataset and evaluate how well a novel view can be interpolated and/or extrapolated with different choices of geometry and textures.

Setup We compare our real-time rasterized renderer (RAS) using geometries reconstructed by the surface based methods of Colmap [48], IDR [60] and our Neural Lumigraph Rendering (NLR). Additionally, we consider both the neural textures generated by the respective method (not available for Colmap) as well as alternative usage of the original training camera images.

We provide the renderer with textures corresponding to 5 of the 6 high-resolution central Back-Bone H7PRO cameras in our dataset (see Section 5), and we measure the PSNR of the held-out interpolated/extrapolated 6th view. The same ground-truth masks are applied for all measurements. We average results from all 6 possible test view choices.

Results Table 4 presents PSNR scores for each scene as well as the overall average. A qualitative comparison is presented in Fig 8. For all choices of textures, we observe a favorable performance of the shape exported from our method compared to shapes exported by both the IDR and Colmap. Colmap generates very accurate but incomplete geometries resulting in holes in the rendered images (note the hair in the scene “L1”). The geometry extracted from IDR provides similar degree of view consistency as our own geometry when combined with the original captured textures, even though our method still maintains a small margin (see the column “Captured textures” in Tab. 4). However, the neural textures generated by IDR lack detail present in our neural textures which results in comparatively lower PSNR scores (see the column “Neural 1 \times ”).

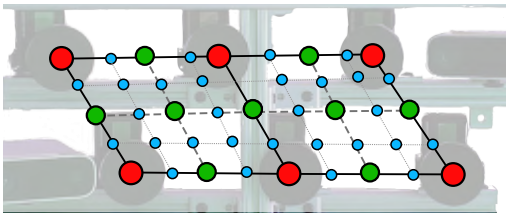


Figure 6. The positions of the original texturing cameras (red) and the interpolation pattern to generate novel textures labeled as “2 \times ” (green) and “3 \times ” (blue).

Scene	Method	Captured textures	Neural textures		
			1 \times	2 \times	3 \times
A1	Colmap	23.96	N/A	N/A	N/A
	IDR	22.23	22.49	23.31	23.61
	NLR (RAS)	23.65	23.83	26.93	30.33
A8	Colmap	11.48	N/A	N/A	N/A
	IDR	19.80	19.65	20.80	21.15
	NLR (RAS)	22.90	23.06	26.00	29.20
L1	Colmap	5.45	N/A	N/A	N/A
	IDR	18.55	18.49	18.89	18.99
	NLR (RAS)	21.24	21.61	24.13	26.56
L3	Colmap	20.03	N/A	N/A	N/A
	IDR	23.56	23.66	25.24	25.95
	NLR (RAS)	24.32	24.85	28.76	32.10
L4	Colmap	13.36	N/A	N/A	N/A
	IDR	17.98	17.73	18.36	18.59
	NLR (RAS)	19.57	19.80	22.20	24.25
M2	Colmap	18.40	N/A	N/A	N/A
	IDR	22.07	21.84	22.60	22.96
	NLR (RAS)	23.94	24.34	26.96	29.76
P4	Colmap	24.22	N/A	N/A	N/A
	IDR	16.06	13.85	14.78	15.59
	NLR (RAS)	24.16	24.27	27.41	30.04
Average	Colmap	16.70	N/A	N/A	N/A
	IDR	20.04	19.67	20.57	20.98
	NLR (RAS)	22.83	23.11	26.06	28.89

Table 4. PSNR scores for the interpolated/extrapolated views achieved by different variants of our real-time renderer. Neural textures for the original camera poses (1 \times), and their supersampled variants (2 \times , 3 \times) are compared to a direct use of the original images captured by our camera array.

Texture space upsampling An important argument for using the neural textures rather than the original captured images is the possibility to use the neural renderer to produce additional virtual views. This allows to generate much denser virtual camera spacing, i.e. through view synthesis, than would be practical with physical camera setups. To demonstrate this unique capability we resampling the texture space by subdividing the original 3×2 camera layout in our dataset as shown in Figure 6. This yields first the original 6 textures (labeled “1 \times ”), then additional 9 textures (total of 15; labeled “2 \times ”) and finally additional 30 textures (total of 45; labeled “3 \times ”). We follow the same procedure and remove the texture corresponding to the original ground-truth camera pose for the test.

The results are again presented in Table 4 (two right-most columns) and in Fig 8. Both IDR and SineSDF experience significant quality boost as the texturing angular space gets denser and the interpolation distances between blended textures get smaller. However, only our technique can leverage the high spatial detail featured in our neural textures and, as a result, shows significantly larger PSNR

overall. Please explore the video supplement to evaluate dynamic properties of this behavior.

B. Additional Results

In this section, we provide additional results that did not fit into the main paper. Table 5 complements Table 2 in the main paper and shows results of the image reconstruction metrics computed for the three test views held-out from the training on the DTU dataset [22]. Table 6 extends Table 3 in the main paper and shows additional image metric scores for results using variety of mutli-view datasets.

C. Supplemental Video

We provide a detailed video supplement to fully evaluate view consistency of our novel rendered views and compare to the baseline methods. A summary video with side by side comparison and comments is available in the root folder as 4945_supplemental_video.mp4 (x264 encoding in mp4 container). Additionally, individual results in form of short videos are provided separately in three folders, each showing a different type of comparison.

Neural renderings Videos in the folder `neural_videos` present neural renderings of scenes from the datasets used in our paper produced by our method as well as well Colmap, IDR, Neural Volumes and NeRF. Note, that the renderings are provided without post-processing and show the background reconstruction present in Neural Volumes and NeRF. This background was not evaluated in metrics in our paper; all tested methods use the same foreground mask when evaluating the PSNR.

Shape renderings In the folder `shape_videos`, we use a simple OpenGL renderer using the PyRender Python package to visualize meshes extracted from Colmap, IDR and our Neural Lumigraph that are later used for the real-time rendering. Note that IDR contains lot of boundary geometry noise for views that were not supervised during training.

Real-time renderings Finally, in the folder `cg_videos`, we present outputs of our real-time renderer as well as different variants presented in the study in Sup. Sec. A.

D. Input sensitivity

Some of the key input properties that affects quality of multi-view reconstruction is the accuracy of camera calibration. We investigate how our method performs if this factor is taken out of the equation by measuring its performance on a computer-generated 3D scene. We have used Unity to

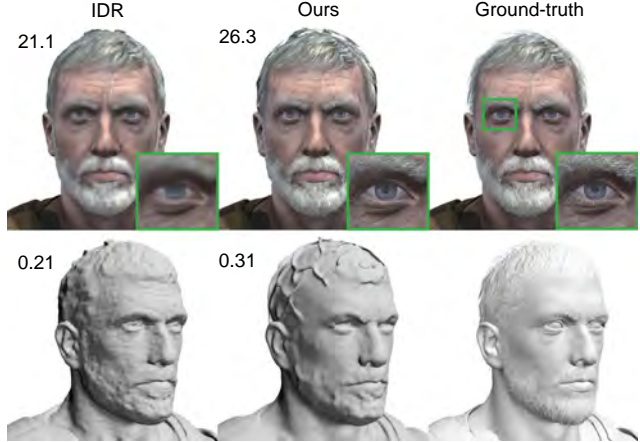


Figure 7. A comparison of image and shape reconstruction by IDR and our method with the ground-truth Unity rendering and mesh shown on the right. The upper numbers denote the image PSNR averaged over all views and the lower numbers correspond to the Chamfer distance.

render a 5×5 grid of views of a 3D head model with surface specular properties at a resolution of 2160×2160 pixel and we extracted the ground-truth camera poses as well as the object masks.

As Fig. 7 shows, both methods achieve a clean shape reconstructions without artifacts. The hair region is represented by a cloud of semi-transparent billboards in the original 3D model, so this would be challenging to be accurately reconstructed by any approach. While the Chamfer distance metric favors the results of IDR, we see that our image quality is much higher as it features notably higher spatial fidelity. We conclude that our method can avoid geometry artifacts if consistent camera and image labeling can be obtained.

E. Datasets

Our Captured Dataset Each capture consists of 22 images and corresponding foreground masks. Six images are captured with high-resolution narrow field-of-view (FOV) cameras, while the remaining 16 images are captured with low-resolution wide FOV cameras. While these additional 16 images provide useful geometrical information, they do not carry high quality visual detail. This is why we use all views for training each of the evaluated methods but we only evaluate reconstruction quality for the 6 central high-resolution cameras.

Figure 16 shows scenes in our dataset. These data are available on the project website.²

²<http://www.computationalimaging.org/publications/nlr/>

Scan	Colmap [48]			IDR [60]			NeRF [37]			NLR-ST			NLR-RAS		
	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓
65	15.12	0.845	0.156	21.89	0.939	0.105	27.08	0.948	0.069	<u>26.58</u>	<u>0.941</u>	<u>0.083</u>	25.60	0.940	0.087
97	11.37	0.835	0.166	22.77	0.914	0.128	24.17	0.918	<u>0.109</u>	26.43	0.922	0.108	<u>26.40</u>	<u>0.922</u>	<u>0.108</u>
106	16.69	0.875	0.153	20.49	0.891	0.205	30.17	0.934	0.107	<u>29.60</u>	0.939	0.098	28.52	<u>0.938</u>	<u>0.104</u>
118	21.72	0.912	0.100	23.24	0.937	0.150	31.03	0.955	0.083	<u>30.77</u>	<u>0.950</u>	<u>0.091</u>	30.39	0.950	0.093

Table 5. Image error metrics PSNR, SSIM [58] and LPIPS [61] computed for the 3 held-out test views of the DTU dataset [22] complementing the respective training errors in Table 2 of the main paper. Best scores in bold, second best underlined.

Dataset	Colmap			IDR			NV			NeRF			NLR-ST			NLR-RAS		
	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓
Volucap (1)	19.57	0.975	0.037	22.27	0.978	0.043	<u>29.32</u>	0.975	0.034	32.71	0.982	<u>0.026</u>	28.59	<u>0.981</u>	0.022	28.57	<u>0.981</u>	0.022
Digit. Ira (1)	Fail.			23.89	0.835	0.286	26.54	0.829	0.287	<u>31.18</u>	<u>0.845</u>	0.267	31.63	0.851	0.255	31.13	0.851	<u>0.260</u>
Ours (7)	17.58	0.841	0.187	22.34	0.878	0.202	25.05	0.857	0.186	28.45	0.898	0.171	30.45	0.921	0.147	<u>30.25</u>	0.921	<u>0.151</u>

Table 6. Average image reconstruction error metrics PSNR, SSIM [58] and LPIPS [61] computed across datasets (number of scenes in parentheses). (*) Only 5 scenes tested.

DTU We use the multiview images and camera calibrations provided by [22] along with object masks by Niemeyer *et al.* [39] and Yariv *et al.* [60]. Each scan consists of 49 or 64 images with 1600×1200 pixels and covering approximately 90×90 degree view zone of the object. Ground truth 3D shapes in form of point captured using structured light stereo are used to compute the Chamfer distance.

Volucap We evaluate a full body reconstruction using a sample frame from a video sequence provided by the volumetric production studio Volucap GmbH [16]. The cameras are distributed in pairs around an upper hemisphere of the capture dome with a subject standing in the center. The images are cropped and rescaled to 2028×1196 px. The masks are manually annotated on top of automatic background subtraction.

Digital Ira The Digital Ira dataset [12] contains 7 very high resolution (3456×5184 px) closeups of a man’s face in a studio setting. We used the provided camera calibration data and we manually annotated objects masks for each view.

F. Baselines

Colmap We follow the same procedure and settings for Colmap [48] as Yariv *et al.* [60] to reconstruct a 3D point cloud from the multi-view images, build a surface mesh using Poisson reconstruction (with trim = 7) and render the images using PyRender. For the DTU dataset, we also apply object masks to remove background points before the surface reconstruction. Note that this method failed to produce a surface mesh of the Digital Ira [12] with all variations of

parameters that we tried.

NeRF We use the original code shared by the authors [38]. We modify the code to support general camera models. We execute the code using the provided high-quality configuration used for their paper results but we reduce the number of random samples to 1024 to fit to our GPU memory (11 GB on Nvidia Geforce RTX 2080Ti). Further, to support variable sensor sizes in our dataset, we upsample all images to the shape of the high-resolution cameras which retains the detailed information.

For our dataset, we are forced by the memory requirements to reduce the training resolution by factor of 2 from 3000×4000 pixels to 1500×2000 pixels. However, it is unlikely that this reduction affected the results since our experiment with even more aggressive four-fold downsizing (to 750×1000 pixels) yielded a PSNR of 28.62 dB which is comparable to 28.58 dB of the two-fold downsizing we ended up using.

IDR We use the original code and configuration shared by the authors [60]. Note, that while IDR supports training without accurate camera calibrations, we train with the same calibrations used for all other methods and keep the poses fixed during the process.

Neural Volumes We use the original code shared by the authors [31]. We train our models for at least 150K iterations. As recommended by the authors, we use three views from different angles to condition the autoencoder. The authors provide the ability to either pass the background explicitly or to let the network estimate the background. Since we don’t have a background image for all datasets,

we choose to let the network estimate the background. Due to the long training times, we evaluate the method on a limited subset of our dataset. We crop the lower resolution landscape images to the same aspect ratio as the six high-resolution central cameras in portrait mode. Due to GPU memory limitations, we are only able to train with the image resolution reduced to 1/4 of the original size of the high-resolution images. Note that Neural Volumes uses an explicit voxel representation and is therefore, unlike implicit models, not resolution independent.

G. Metrics

PSNR For all methods we compute the PSNR between the reconstructed images and the ground truth only for pixels in the object masks. The same procedure is used by Yariv *et al.* [60] and we verify it by reproducing their metric scores with their pre-trained models. Therefore, even methods that do not use masks to recover the shape (NeRF, Neural Volumes) are only rated based on prediction of the foreground object.

Chamfer distance The Chamfer distance is computed in a similar manner as in the original Matlab scripts provided with the DTU dataset [22]. The shortest distance between each ground-truth 3D point and the reconstructed shape is computed in one direction by the Open3D [63] library for Python. This prevents penalizing reconstruction of background (NeRF) or reconstruction of occluded object parts. We leave out Neural Volumes from evaluating the metric since surface extraction was not demonstrated in the original manuscript [31] and the explicit nature of the representation makes recovery of fine details difficult.

H. Sphere tracing implementation details

Convergence of the main paper Eq. 5 towards the zero-level set S_0 is not guaranteed for general shape configuration. As the SDF function S indicated distance towards the nearest surface, it underestimates the optimal step length if the nearest surface is not orthogonal to the current ray. This can partially be improved by dividing the step length by the dot product $-\mathbf{r}_d \cdot \nabla_{\mathbf{x}} S(\mathbf{x}_i)$. However, computing the gradients in every step is costly and quickly diminishes the returns. Furthermore, if a ray closely misses a surface edge with a surface normal orthogonal to the ray direction, such adjustment does not improve the convergence as the step length will stay very small.

To mitigate these issues, we use bidirectional sphere tracing as described by Yariv *et al.* [60]. This approach does not rely on the sphere tracer’s convergence as it uses it to only narrow a possible range of surface positions which is then further refined by additional solvers.

To this goal, we solve Eq. 5 to get the near zero-level set \mathbf{x}_n^n in a forward direction for $n = 16$. We mask rays with $|S(\mathbf{x}_n)| < 5e^{-5}$ as converged. These rays are not further optimized and $\hat{\mathbf{x}}_n = \mathbf{x}_n^n$ is the final output of the sphere tracer. Next, for the remaining rays, we solve a modified equation in an opposite direction to find a far zero-level set \mathbf{x}_n^f as

$$\mathbf{x}_0^f = \mathbf{r}_o + t_f \mathbf{r}_d, \quad \mathbf{x}_{i+1}^f = \mathbf{x}_i' - S(\mathbf{x}_i^f) \mathbf{r}_d. \quad (14)$$

where t_f is the rear intersection of \mathbf{r}_d with a unit sphere. \mathbf{x}_n^n and \mathbf{x}_n^f define the nearest and the furthers possible location of the first surface point \mathbf{x}_n . In cases where \mathbf{x}_n^f is closer than \mathbf{x}_n^n , we conclude that no surface exists. In the remaining cases we evaluate 100 samples with the candidate range and look for the first zero crossing of the S . Finally, assuming a locally monotonic S we further refine the zero-crossing location by 8 steps of sectioning [60] to get the zero-level set $\hat{\mathbf{x}}_n$.

Note that for performance, stability and memory efficiency, this procedure is not differentiated [60, 24, 30]. Therefore, one extra step of the forward sphere tracing is executed with gradient tracking enabled to achieve a differentiable sphere tracing. We apply the gradient direction adjustment [60] for this last step to get the final zero-level set as

$$\mathbf{x}_n = \hat{\mathbf{x}}_n - \frac{S(\hat{\mathbf{x}}_n)}{\nabla_{\mathbf{x}} S(\hat{\mathbf{x}}_n) \cdot \mathbf{r}_d} \mathbf{r}_d, \quad (15)$$

for points where $|S(\mathbf{x}_n)| < 0.005$. A division by zero is avoided by clamping the denominator to $|\nabla_{\mathbf{x}} S(\hat{\mathbf{x}}_n) \cdot \mathbf{r}_d| > 0.01$.



Figure 8. Comparison of our real-time rendering variant evaluated on an view not present in the texture set. The IDR results for the presented view in P4 are black due to an occlusion by a false geometry generated in front of the person (see the shape rendering). While fixable and only present in this scene, we keep the geometry as is and provide detailed metrics for each single scene in Table 4. A rejection of this result would not change the overall trend.



Figure 9. A1



Figure 10. A8



Figure 11. L1

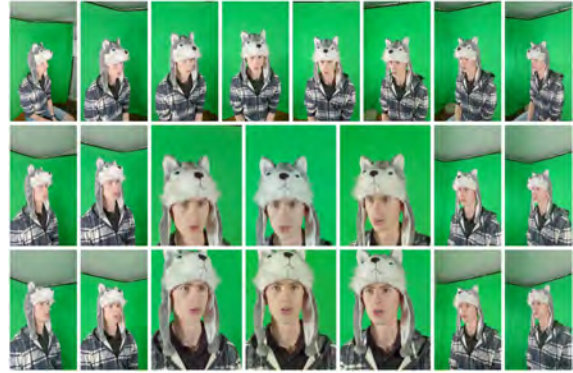


Figure 12. L3



Figure 13. L4



Figure 14. M2



Figure 15. P4

Figure 16. The seven scenes forming our dataset. Image positions approximate the camera layout. The central 6 images are high-resolution head close-ups.