

Towards Learning Neural Representations from Shadows

Kushagra Tiwary ^{*}, Tzofi Klinghoffer ^{*}, and Ramesh Raskar

Massachusetts Institute of Technology
{ktiwary, tzofi, raskar}@mit.edu

Abstract. We present a method that learns neural shadow fields which are neural scene representations that are *only* learnt from the shadows present in the scene. While traditional shape-from-shadow (SfS) algorithms reconstruct geometry from shadows, they assume a fixed scanning setup and fail to generalize to complex scenes. Neural rendering algorithms, on the other hand, rely on photometric consistency between RGB images, but largely ignore physical cues such as shadows, which have been shown to provide valuable information about the scene. We observe that shadows are a powerful cue that can constrain neural scene representations to *learn* SfS, and even outperform NeRF to reconstruct otherwise hidden geometry. We propose a graphics-inspired differentiable approach to render accurate shadows with volumetric rendering, predicting a shadow map that can be compared to the ground truth shadow. Even with just binary shadow maps, we show that neural rendering can localize the object and estimate coarse geometry. Our approach reveals that sparse cues in images can be used to estimate geometry using differentiable volumetric rendering. Moreover, our framework is highly generalizable and can work alongside existing 3D reconstruction techniques that otherwise only use photometric consistency.

Keywords: Scene Representations, Differentiable Rendering, 3D Scene Reconstruction, Shape-from-Shadows, Volume Rendering

1 Introduction

Recovering 3D geometry from 2D images remains an extremely important, yet unsolved problem in computer vision and inverse graphics. Considerable progress has been made in the field when assumptions are made, such as bounded scenes, diffuse surfaces, and specific materials. However, reconstruction algorithms still remain largely susceptible to real world effects, such as specularities, shadows, and occlusions [34]. This susceptibility is largely due the variation in different materials and textures, and a non-unique mapping from 3D geometries to 2D images. Even though these effects cause issues for many methods, they also provide valuable information about the scene and geometry of the object. For example, cues like self-shadows provide vital information about an object’s concavities, while shadows cast on the ground plane provide information about its geometry. Moreover, shadows are independent of textures and surface reflectance models and

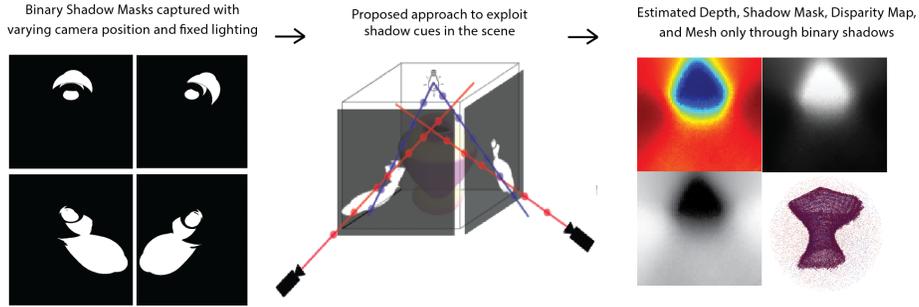


Fig. 1. Exploiting physical cues in neural rendering. Our approach takes sparse binary shadow masks captured with varying camera positions under fixed lighting and uses our proposed differentiable shadow rendering model to estimate shadow maps, thereby learning neural scene representations. We can visualize the learned implicit representations by rendering estimated depth maps and estimated shadow maps from novel views. We also run marching cubes [15] on our learned representations to get explicit meshes for a quantitative analysis.

are a strong cue in overhead imagery where vertical surfaces, like facades, are sampled poorly, whereas oblique lighting can expose this geometry. Exploiting, instead of ignoring these cues, can make algorithms robust and the fundamental problem of 3D reconstruction less ill-posed.

Previous works in recovering 3D shape of objects by exploiting physical cues has relied on constructing inverse models to explicitly handle and exploit cues such as shadows, shading, motion, or polarization [2] [41] [40]. These approaches are physically anchored as they use properties of light or surface reflectance models to exploit cues and only need up to a single image to reconstruct simple objects. Albeit successful under strict assumptions about lighting, camera, and the object, these models typically cannot handle complex scenes and do not translate well into real-world scenarios as creating inverse models to capture complex physical phenomenon soon becomes intractable and hard to optimize.

To combat the problem of real world variability, modern methods such as [31] [17] [22] [28] [37] [13] have largely been data-driven by directly learning 3D representations on real-world scenes based on photometric consistency. Such methods employ an *analysis-by-synthesis* approach to solve the problem by using machine learning to search the space of possible 3D geometries and an inverse model to synthesize the scene based on the predicted geometries. These approaches typically only optimize the photometric loss between different camera viewpoints and show success in learning implicit representation by rendering novel views. However, because they do not explicitly handle these physical cues in their forward model, they fail in scenarios with complex lighting [29], specularity [39], or reflections [6].

* Equal contribution

Motivated by the above observations, we explore what can be learned by exploiting physical cues in a data-driven neural rendering framework. In this paper, we investigate whether the neural rendering framework can learn geometry from physical cues without the assumptions made by the aforementioned methods. We study the use of shadows cast by objects onto themselves and nearby surfaces as the only source of information for 3D reconstruction. While modern approaches for 3D reconstruction ignore such cues, we aim to exploit them. Our unsupervised approach uses *only* shadows to reconstruct the scene by leveraging recent advances in volumetric rendering and machine learning, and therefore proposes a physically anchored data-driven framework to the problem of shape from shadows. Moreover, unlike previous work in shape from shadows, we present a novel method that uses differentiable rendering in the loop to iteratively reconstruct the object based on a loss function instead of iteratively refining the object through explicit carving. Specifically, we use an efficient shadow rendering technique called shadow mapping as the forward model and make it differentiable so that it can be used as an inverse model to iteratively reconstruct the object. Our work also reveals that from limited cues the differentiable volumetric rendering component can *quickly converge to localize and reconstruct a coarse estimate of the object when such cues are explicitly modeled by a forward model*. Our work also suggests that neural rendering can exploit shadows to recover hidden geometry, which otherwise may not be discovered by photometric cues.

1.1 Contributions

Our contributions in this paper are the following:

- A framework that directly exploits physical cues like shadows in neural renderers to recover scene geometry.
- A novel technique that integrates volumetric rendering with a graphics-inspired forward model to render shadows in an end-to-end differentiable manner.
- Results showing that our framework can learn coarse scene representations from just shadows masks. We evaluate the learned representations qualitatively and quantitatively against vanilla neural rendering approaches. To the best of our knowledge, we are the first to show that it is possible to learn neural scene representations from binary shadow masks.

2 Related Work

Shape from Shadows. Shadowgram imaging deals with estimating the shape of an object through a sequence of shadow masks captured with light sources at various locations. These methods typically assume a controlled and fixed object scanning setup [27] [36]. Martin & Aggarwal [16] introduced a volumetric space carving approach to SfS which outputs a visual hull around the object by carving out voxels lying outside the visual cone. Other work takes a more probabilistic

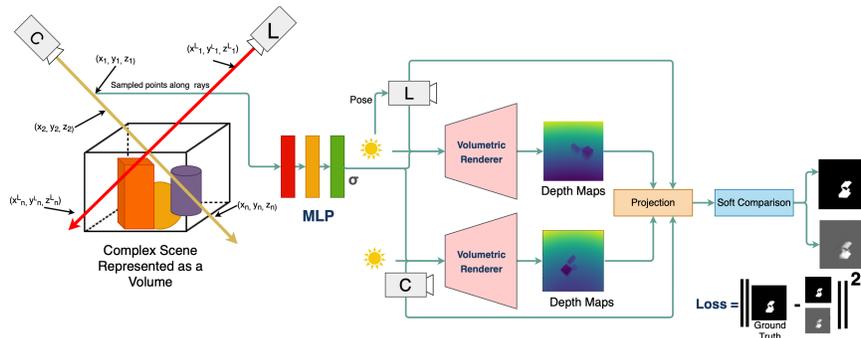


Fig. 2. Overview of the proposed pipeline We train a neural network to predict opacity at points along the camera and light rays. The opacities are used by the volumetric renderer to output the ray-termination distance which we use to estimate the *z-buffer* from the camera and the light perspective, the latter also known as the shadow map. The estimated *z-buffer* is fed into a **Projection** step that projects the camera pixels and their associated depths into the light’s reference frame. The shadow map is indexed to obtain the corresponding depth values at these new points. The projected depths and indexed depths go through a **Soft Comparison** step which outputs predicted cast shadows in the scene from the camera’s perspective. A loss is computed on the *predicted* and the *ground-truth* shadow mask.

approach to the shape-from-silhouettes problem to make the algorithm more robust to errors [9]. However, interpreting shadows as silhouettes means that self-shadows are not handled, thus motivating Savarese et al. [27] to propose a method to “carve” out objects based on self-shadows to create more complete reconstructions.

In contrast, our work takes a differentiable approach to solving the problem through learning. Instead of an explicit carving of voxels we first construct a differentiable forward model that casts shadows based on some geometry. Then, we let the machine learning component predict geometry, which is synthesized by the renderer to cast shadows. Finally, we optimize this setup based on a mean square error between predicted and ground truth shadow masks.

Neural Rendering Broadly speaking, a neural rendering framework is composed of a differentiable renderer, which can render the scene based on input parameters and is able to differentiate the scene w.r.t. those input parameters. While there are many formulations of differentiable renderers [21] [14] [10] [8] that can synthesize scenes, state-of-art approaches have shown tremendous success by relying on differentiable volumetric rendering [20]. Volumetric rendering approaches can realistically render complex scenes and are gradient-friendly. Thus, typical approaches train a neural network to encode the scene and optimize it for photometric consistency between input 2D images from different viewpoints [17] [28] [18] [19]. Recent methods such as [6] [33] [29] [3] explicitly account for specular, reflections and other such phenomenon, however, the

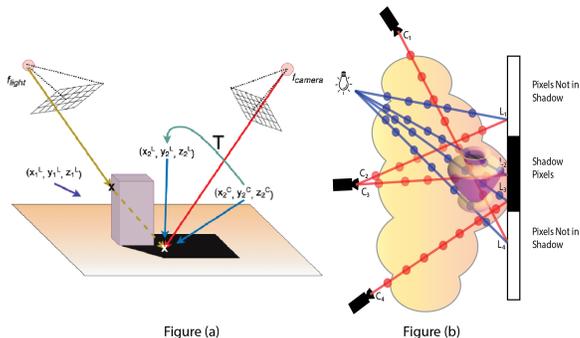


Fig. 3. Figure (a): A point $\mathbf{x} \in \mathbb{R}^3$ in the scene is defined to be in shadow if no direct path exists from the point \mathbf{x} to the light source, implying that there **must** be an occluding surface between \mathbf{x} and the light source. We differentially render the scene’s depth from the camera and the light’s perspective at each pixel and then project the camera pixel and its depth into the light’s frame of reference. We then index the light’s depth map, or z-buffer, to get z_1^L . We note that z_1^L is less than z_2^L , i.e. there must be an occluding surface as a ray projected from the light’s perspective terminates early. This implies that this point is in shadow. Figure (b) shows a 2D slice of our approach and represents a volume (cloud) with the shadow mask unraveled. The network learns an opacity per point (dots) via the shadow mapping objective which penalizes predicted geometries that don’t cast perfect ground truth shadows. Through this, the networks learn 3D geometry that is consistent across all shadows maps for all cameras given a particular light source.

goal of these works are to improve novel view synthesis. Thus, these methods still rely on learning the scene using photometric information.

In contrast, our work deals with 3D reconstruction, not novel view synthesis, and explores what can be learned by relying on shadow cues in the scene. Our framework only operates in the shadow input and output space to infer a 3D representation of the scene. In addition, similar to [26] [38], our work also reveals that differentiable volumetric rendering is a powerful component that can learn the scene by only relying on sparse physical cues. While volumetric approaches rely on a photometric cues, differentiable rasterization [8] [12] has been shown to reconstruct 3D mesh using single low dimensional images of ShapeNet objects [4] by only using silhouettes. However, these methods fail to show success on high dimensional images, while our approach can scale up to higher dimensional images. Concurrent work by Liu *et al.* [11] also leverages shadows to perform 3D reconstruction, but integrates learned object priors, whereas we solely rely on binary shadow masks and use volumetric rendering.

Shadows in Graphics. Graphics deals with the forward model and shadow mapping [35] is one of the most efficient techniques to render shadows in a scene given the scene’s geometry, camera viewpoint and light position. While differentiability is not important for graphics, we make the shadow mapping

framework differentiable to work with modern 3D reconstruction algorithms. We describe the algorithm and our implementation in Section 3.

3 Neural Representations From Shadows

Our goal is to recover the scene through shadows cast on the other objects or onto itself. Our method recovers shadows in an image by applying a threshold on that image thereby making no distinction between types of shadows. We show how we model the shape-from-shadows problem using differentiable rendering and implicit representations in Section 3.1 and our graphics-inspired differentiable forward model in Section 3.2. In Section 3.3, we discuss our additional techniques that we use to enable optimization on binary shadow masks.

3.1 Scenes as Neural Shadow Fields

Implicit Scene Representations. Similar to Mildenhall *et al.* [17], we represent a continuous scene by parametrizing it using a learnable function f_θ . However, our approach does not include any photometric component, therefore we represent the scene as a 3D function with input $\mathbf{x} = (x, y, z)$ and a volumetric density σ as output.

$$\gamma(\mathbf{x}) = \left(\sin(2^0 \pi \mathbf{x}), \cos(2^1 \pi \mathbf{x}), \dots, \sin(2^L - 1 \pi \mathbf{x}), \cos(2^L - 1 \pi \mathbf{x}) \right) \quad (1)$$

$$f_\theta : \mathbb{R}^L \rightarrow \mathbb{R}^+; (\gamma(\mathbf{x})) \mapsto (\sigma)$$

We use a positional-encoded 3D point $\gamma(\mathbf{x}), \{\gamma(\mathbf{x}) \in \mathbb{R}^L, \mathbf{x} \in \mathbb{R}^3\}$ as input, which maps to an associated volumetric density $\sigma \in \mathbb{R}^+$ [17] [30]. In contrast, f does not encode view dependant color and is independent to viewing direction.

Volumetric Renderer. We define a volumetric renderer \mathbf{R}_{vol} that takes N opacities $\{\sigma\}_{i=1}^N$ at N discretely sampled points $\{\mathbf{x}\}_{i=1}^N$ along a ray \mathbf{r} .

$$\mathbf{R}_{\text{vol}} : [\mathbb{R}^+]_{i=1}^N \rightarrow [\mathbb{R}^+]_{i=1}^N; (\{\sigma\}_{i=1}^N) \mapsto (\mathbf{d}) \quad (2)$$

Since we only have binary shadows as input, we modify the renderer to output the ray termination distance, \mathbf{d} , instead of the radiance at that ray. \mathbf{R}_{vol} is not a trainable component, but the ray termination distance, \mathbf{d} , is differentiable w.r.t. the input opacities. The estimated ray-termination distance, range, is computed as follows:

$$\hat{\mathbf{D}}(\mathbf{r}) = \sum_{i=1}^N T_i \alpha_i t_i; T_i = \prod_{j=1}^{i-1} (1 - \alpha_j); \alpha_i = (1 - e^{-\sigma_i \delta_i}) \quad (3)$$

We sample $\mathbf{r}(t)$ at points $\{t_0, \dots, t_N\}$ and evaluate the function $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ to get sampled points $\{x_0, \dots, x_N\}$ in the scene. T_i is defined as the cumulative transmittance from t_0 until t_i and $\delta_i = t_{i+1} - t_i$ which is the distance between two samples. σ_i is the estimated opacity at point i by a learned function f_θ . Intuitively, the renderer gives us the ray termination distance for each ray shooting through a pixel.

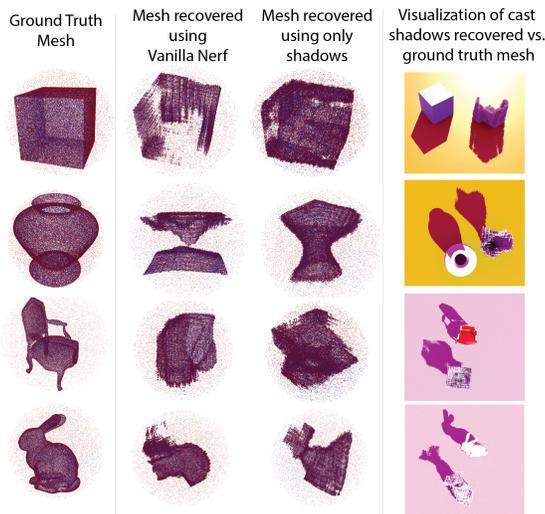


Fig. 4. Qualitative Results. We observe that for overhead views of the scene where the vertical surface of the vase is sampled poorly in the RGB space, vanilla NeRF fails to exploit geometry cues hidden in cast shadows compared to our approach. Our method doesn’t impose any object priors therefore it infers a geometry that will minimize the difference between the predicted and true shadow. Column 4 illustrates that rendered shadows are very similar, indicating that the differentiable rendering framework can indeed learn geometry from sparse shadow cues. Some parts of the objects such as the upper face of cuboid are never in shadow, therefore our approach yields no reconstruction for those surfaces, further showing that the geometry is indeed *only* learnt from cast shadows. We extract the mesh from the volume using marching cubes and visualize it here using a point-cloud SDF representation.

3.2 Differentiable Shadow Mapping

We define any point $\mathbf{x} \in \mathbb{R}^3$ in the scene to be in shadow if no direct path exists from point \mathbf{x} to the light source \mathbf{L} . This logic implies that there **must** be some object or an occluding surface between the point \mathbf{x} and \mathbf{L} that occludes the light ray from reaching point \mathbf{x} . In graphics, shadow mapping [35] uses this observation to construct a forward model to render efficient and accurate shadows in the scene based on known light and camera sources. Our approach makes this efficient shadow rendering forward model differentiable so that it can be used as an inverse model. We then pose the problem of shape from shadows and use our proposed inverse model to estimate the 3D geometry of the scene.

Estimated z-buffer. We first evaluate the renderer from the camera’s perspective to get the estimated ray termination distance, or range map, $\hat{\mathbf{D}}_{cam}$ for all rays coming out of the binary shadow map. However, shadow mapping requires the depth perpendicular to the image plane, i.e along the z axis of the camera’s local coordinate system. This depth is equivalent to a *z-buffer* in graphics and we refer to this value as the *depth* at that pixel. We define a function g to estimate

the *z-buffer* $\hat{\mathbf{Z}}$ from the range map $\hat{\mathbf{D}}$.

$$\hat{\mathbf{z}}_{u,v} = g(\mathbf{d}_{u,v}) = \frac{\mathbf{d}_{u,v}}{\|(u, v, 1) \cdot \mathcal{E}\|_2} \quad (4)$$

The function takes a ray shooting from a pixel (u, v) and a predicted range, $\hat{\mathbf{D}}_{cam}^{u,v}$ as input. \mathcal{E} is the rotational component of the camera’s extrinsic matrix, $\mathbf{d}_{u,v}$ is the ray termination distance from camera’s focal point, and $\hat{\mathbf{z}}_{u,v}$ is the depth along the *z*-axis from the pixel (u, v) . We also compute the estimated *z-buffer* from the light’s perspective, which we refer to as the estimated *shadow map*.

Projection. With the estimated depths at each pixel from the camera and the light source, we now need to estimate which camera pixels are in shadow given the particular light source. As illustrated in Figure 3, we do this by projecting all pixels and their associated depths visible by the camera into the light’s frame of reference. We then use this projected coordinate to index the shadow map to get the depth to that point from the light’s perspective. We formally write this as follows:

$$\begin{aligned} (U_{cam}^l, V_{cam}^l, \hat{\mathbf{Z}}_{cam}^l) &= (U_{cam}, V_{cam}, \hat{\mathbf{Z}}_{cam}) \cdot P_{light_from_cam} \\ \hat{\mathbf{Z}}_{light}^{U_c^l, V_c^l} &= \hat{\mathbf{Z}}_{light} \left[U_{cam}^l, V_{cam}^l \right] \end{aligned} \quad (5)$$

Here, $\hat{\mathbf{Z}}_{cam} \in \mathbb{R}^{H \times W}$ is the estimated *z-buffer* from the camera’s perspective at pixels $\{U_{cam}, V_{cam}\} \in \mathbb{R}^{H \times W}$. $P_{light_from_cam}$ is the projection matrix to the light’s reference frame from the camera’s. We denote $(U_{cam}^l, V_{cam}^l, \hat{\mathbf{Z}}_{cam}^l)$ as the pixels and depth in camera’s frame (subscript) projected into the light’s frame, denoted by the superscript l . We index the shadow map, $\hat{\mathbf{Z}}_{light} \in \mathbb{R}^{H \times W}$, at the projected camera pixels to retrieve the depth of the projected camera pixels from the light source. This is denoted as $\hat{\mathbf{Z}}_{light}^{U_c^l, V_c^l}$ which is the shadow map indexed at pixel locations U_c^l, V_c^l . In practice, not all pixels will project within the shadow map’s height and width constraints specified at the start of training. In graphics, these pixels are usually ignored, however, we clamp all our projections to lie within the height and width bounds to maintain differentiability.

Soft Comparison. Once we have the depths to the projected camera pixels and the depths from the light source to those pixels in the same reference frame, we can then compare them to discover if the camera pixel is in shadow. As illustrated by Figure 3, if the depth from the light source to a point is less than the depth from the camera projected into the light’s frame, it means that the light ray must have intersected an object before reaching that point. Thus, that point must be in shadow. Based on this logic, we formulate a soft comparison, which compares different depths to output the predicted binary shadow mask as follows:

$$\begin{aligned} \Delta \hat{\mathbf{Z}}_{light} &= \left(\hat{\mathbf{Z}}_{cam}^l - \hat{\mathbf{Z}}_{light}^{U_c^l, V_c^l} \right) \\ \hat{\mathbf{M}}_{binary} &= \mathbf{max} \left(\frac{\Delta \hat{\mathbf{Z}}_{light}}{\beta}, \epsilon \right) \end{aligned} \quad (6)$$

We denote $\hat{\mathbf{M}} \in \mathbb{R}^{H \times W}$ as the output of the entire pipeline: predicted shadow masks. The input to our soft comparison is the projected camera z-buffer into the light’s frame, $\hat{\mathbf{Z}}_{cam}^l$, and the shadow map indexed at the projected points $\hat{\mathbf{Z}}_{light}^{U^l, V^l}$ from the **Projection** step. β is a scaling hyper-parameter used to enlarge or decrease the difference, and ϵ is a threshold. We also formulate a “smoother” version of the predicted shadows:

$$\hat{\mathbf{M}}_{smooth} = \mathbf{S}(\text{normalize}(\Delta \hat{\mathbf{Z}}_{light}, \mu_{min}, \mu_{max})) \quad (7)$$

Here, μ_{min}, μ_{max} are used to control the normalization function and \mathbf{S} is the sigmoid function.

3.3 Optimization

To enable convergence, we smooth the binary ground truth shadow masks \mathbf{M} to better guide the framework in predicting accurate shadow masks.

Distance Transform. Binary images contain limited information for differentiation as the gradient is zero everywhere except for the edges where it is one. To encourage our model to estimate better shadow masks, thereby learning a better 3D model, we use a distance transform on the ground truth shadow masks. Specifically, we scale pixel intensities of a binary shadow mask by their distance to the nearest shadow edge. We modify the weighted distance transform in [25] for our approach. The transformed binary shadow mask, $w(\mathbf{M}, \sigma) = \mathbf{M}_w$ is computed as follows:

$$w(\mathbf{M}, \sigma) = \mathbf{M} + \left(w_c(\mathbf{M}) + w_0 \cdot \exp\left(-\frac{(d_1(\mathbf{M}) + d_2(\mathbf{M}))^2}{2\sigma^2}\right) \right) \quad (8)$$

Here, \mathbf{M} is the ground truth binary shadow mask computed after applying a fixed threshold on binary images. w_c is weight map to balance class frequencies, w_0 and σ are hyper parameters. d_1 and d_2 are distances to the nearest and second nearest cell, respectively. We note from our experiments that this particular distance transform yields the most consistent convergence compared to other distance transforms, such as blurring.

Shadow Mapping Loss. We optimize our entire framework on binary shadow masks and train the MLP on the following loss:

$$\mathcal{L}_{sm} = \|w(\mathbf{M}, \sigma) - \hat{\mathbf{M}}\|^2 \quad (9)$$

Here, $w(\mathbf{M}, \sigma)$ is the σ weighted ground truth shadow mask, and $\hat{\mathbf{M}}$ is the predicted shadow mask from Equation (7).

4 Implementation

4.1 Dataset

We create a dataset of objects, including a cuboid, vase, chair and bunny in blender and render images of size 800×800 . Although our approach does not

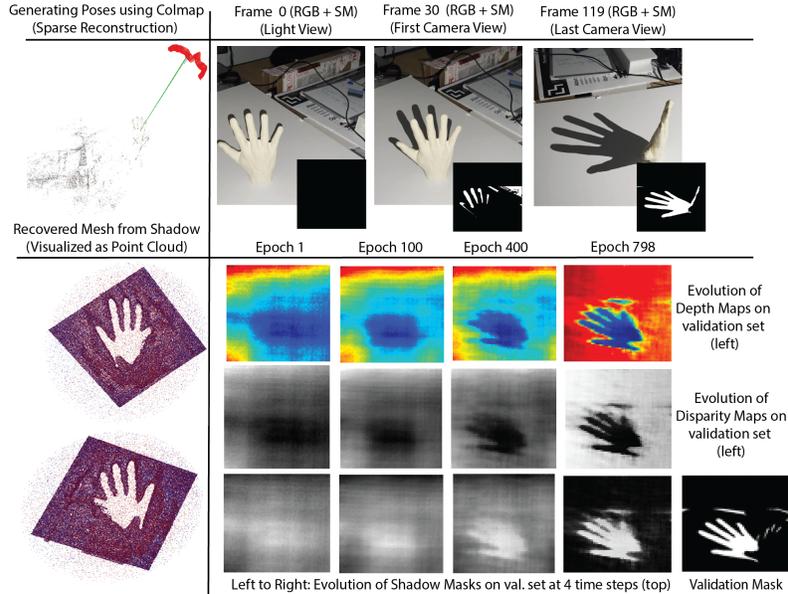


Fig. 5. Real-World Experimentation: We use the *exact same pipeline and training scheme* to reconstruct a 3D mesh from real-world data. We take a video on the iPhone to generate poses for light and camera using COLMAP [43](video link) and extract shadows using an intensity threshold. We show that our method can reconstruct a finer mesh of the hand from the real-world images. We highlight that our method can more easily generalize from sim2real in comparison to photometric approaches since we learn from only shadow masks, which are invariant to many real-world effects, such as texture.

have any constraints on the number or positions of light and camera, we fix one light source and randomly sample 200 camera positions on along the upper half of a sphere around the object. In simulation, we only consider top down/satellite views with the light source being very far away from the camera and the object to represent a distant “sun-like” source although this assumption was relaxed in real world dataset. Our dataset motivates the use of shadows for 3D reconstruction because in overhead imagery, vertical surfaces, like facades, are sampled poorly, whereas oblique lighting can expose this geometry. Link to the dataset is provided in the supplementary section. Details about the real world dataset is also provided in the supplementary section.

4.2 Training Details

We use a faster implementation of NeRF [17] from [24], which uses PyTorch Lightning as its backend [23] [5]. We down sample all images from 800×800 to 64×64 to fit on one RTX-3080 GPU. We use the same positional encoding scheme, $\gamma(x)$ and the MLP configuration f_θ used in [17]. For camera projec-

Scene	RMSE Shadow Mesh	RMSE Vanilla NeRF
Cuboid	0.0078	0.097
Vase	0.010	0.0.011
Bunny	0.0109	0.0106
Chair	0.0092	0.0096

Table 1. We quantitatively analyze the quality of the reconstructed meshes by running ICP [1] on meshes generated by our proposed method, which only uses binary shadows masks, and meshes generated by a vanilla NeRF trained on full RGB images. We show RGB images from Vanilla NeRF in the supplementary along with training details.

tions, we write a custom *Planar Projection Camera* class that encapsulates the projections and readily works with OpenGL and blender cameras. We gradually decrease the σ (Equation 8) from $\{150, 100, 50\}$ during training to encourage the network to learn coarse to fine geometry. We also train a Vanilla NeRF model on RGB images of the same scene on resolution 64×64 . To reconstruct the meshes, we run marching cubes on the learned implicit representations. More information on the exact training details is given in the supplementary section, including details about our more efficient differentiable shadow mapping implementation, which decreases the training time by half.

5 Results

Evaluation Details. We evaluate the performance of our method using root mean square error (RMSE) between the predicted point cloud and the ground truth point cloud, acquired with the iterative closest point (ICP) [1] algorithm, as reported in Table 1. In addition, we assess the visual quality of the predicted depth and shadow masks, and surface mesh, as shown in Figures 4 and 6. The thresholds used to get the mesh from a volumetric representation are given in supplementary material.

Simulated 3D Reconstruction Results. We show the learned scene representations qualitatively by converting them to explicit meshes and rendering them using a signed distance function (SDF). Figure 4 shows the estimated meshes from our method on four object types. We compare our meshes to meshes generated by running vanilla NeRF on RGB images, and the ground truth by running marching cubes on the volume. Our datasets are rendered with overhead camera viewpoints, which enables shadows to be exploited. Given the binary and sparse nature of shadow masks in terms of their information content, we observe that our forward model coupled with the differentiable rendering framework converges to good coarse estimates of object geometry. Moreover, in the case of vases, the mesh reconstruction benefits from exploiting shadows as the algorithm can use *hidden* cues present in the scene, such as the curvature of the vase, which are only partially visible when relying on photometric cues. We also show predicted depth maps and shadow masks on novel camera viewpoints not used during training in the supplementary materials.

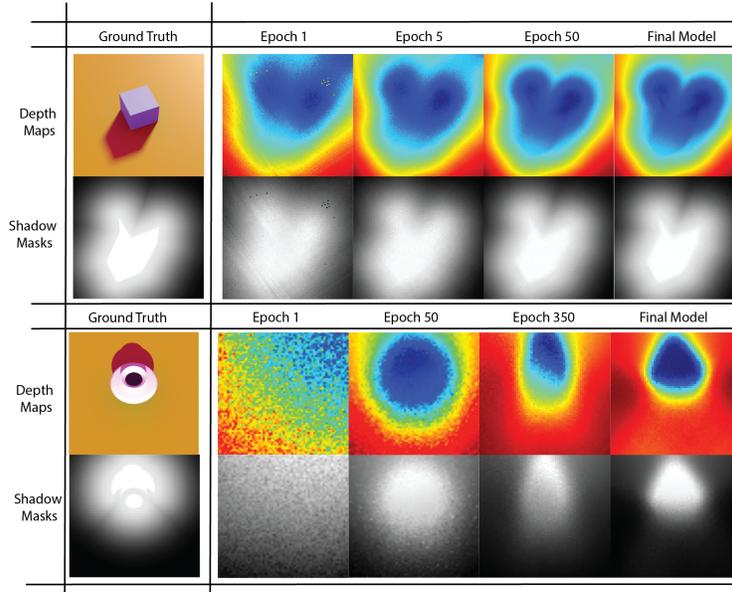


Fig. 6. Evolution of depth maps during training through a novel camera viewpoint. We visualize how the proposed forward model and differentiable rendering framework quickly converges to localize the object based on sparse shadow cues present and then slowly refines the coarse estimate on novel camera viewpoints. We believe that these results reveal that the differentiable volumetric rendering is a powerful framework that can rely on exploiting such physical cues to infer scene information.

Real-World Reconstruction Results. We show our method’s ability to converge to a fine mesh on real-world data of a hand in Fig. 5. Information on data acquisition is provided in the supplementary materials. We first note that our method is robust to coarse light poses as there are visible shadows from the estimated light’s pose in Fig. 5 (please refer to the supplement for details). Our method is able to converge to a fine mesh of the hand, including the fingers and the space between them. We use only 74 shadow masks which makes our method versatile to environments with limited camera views and rarer objects, and no object priors. Moreover, we also show the convergence of the estimated shadow masks, disparity and depth maps from a novel viewpoint. The final estimated shadow mask shown in Fig. 5 is similar to the validation shadow mask and also contains some shadow artifacts due to the threshold segmentation. Additionally, the sim2real gap is not present for our pipeline as it only uses object shadows.

Lastly, we briefly discuss how the data-driven components finds the easiest solution that is consistent with our physics-driven forward model but not the actual world. We note that our training data only has views of cast shadows and does not contain any self-shadows which are present on the back of the hand. This causes the algorithm to instead estimate the mesh of the table and create

a hollow imprint of the hand such that the specified shadow constraint is met. We note that a stand-alone mesh of the hand can be recovered from this imprint and that the recovered mesh is a possible solution given the shadow masks and proposed model. Imposing object priors or adding an extra view of self-shadow to the training set could result in a stand-alone mesh.

Novel Viewpoint Rendering. We observe predicted depth and shadow masks rendered from novel viewpoints in Figure 6. The depth maps converge quite quickly to localize the object even when optimizing on the sparse physical cue of shadows. We posit that this convergence shows how powerful differentiable rendering is for exploiting physical cues to enable better 3D reconstruction. The depth maps converge slowly and we nudge the convergence by gradually decreasing the sigma values for the distance transform. The use of the distance transform leads to blurrier boundaries, however, the rendered mesh shows that a reasonably coarse 3D estimate is captured.

Quantitative Analysis. We also run our datasets on a vanilla NeRF [17] implementation [24]. At lower resolutions and overhead viewpoints, we see that the NeRF approach fails to provide a reasonable fine mesh. We believe this failure is due to the down-sampling of images to 64×64 , which may also be a reason as to why our meshes fail to capture fine details. We run ICP [1] on the generated points cloud and show on-par results to the NeRF approach. Our goal, however, is not to outperform NeRF but to show the effectiveness of differentiable rendering framework in exploiting physical cues instead of ignoring them. The main takeaway from Table 1 is that differentiable volumetric renderers do not need to rely on 8 bit RGB information to reconstruct accurate meshes, but can also leverage other sources of information in the image in addition to relying on photometric cues.

Limitations. In cases such as the cuboid and the vase, we observe that the renderer converges to a predicted mesh that minimizes the shadow masks and the predicted shape even though it is typically a coarse estimate that envelopes the entirety of the object. This means that we see artifacts such as the pointed curve in the vase mesh, or the curvature of the bunny. Since our algorithm only has geometry information where the binary shadow mask is true, areas that are never in shadow have no surface, which leads to incomplete meshes. Imposing a prior can be a solution to this problem. Moreover, our method also assumes known lighting position, which may not always be available.

6 Discussion

Exploiting Physical Cues. One of the major goals of our work is to propose a framework within neural rendering that can readily exploit and learn from, instead of ignore, sparse physical cues, such as shadows. We believe that Fig. 5 shows that sparse physical cues like shadows, actually encode a lot of *hidden* information about the scene and can indeed be exploited. By constructing explicit differentiable forward models and leveraging gradient-friendly volumetric

rendering, we can exploit these cues in conjunction with relying on photometric consistency between images.

Differentiable Shadow Rendering. In rasterization, shadow computation is done through shadow mapping as it is well suited and efficient. However, shadow computation in ray tracing are expensive as every ray needs to compute a path to the light source. Therefore, many ray tracing approaches also use shadow mapping to compute shadows efficiently. We use shadow mapping in our neural rendering approach as well. Our approach is similar to the shadow mapping in graphics as it assumes a binary label on shadows and does not consider soft shadows or ambient lighting. However, we invert shadow mapping and exploit it to do 3D reconstruction, not to render photorealistic images. Moreover, our approach is readily extendable to varying light and camera sources.

6.1 Future Work

We observe that volumetric rendering can converge onto coarse estimates of the object geometry by only relying on shadows, and can be extended to problems such as non-line-of-sight imaging (NLOS) [32] and imaging behind occluders [7]. As shadows themselves are never the only cue present to reconstruct the scene, our work can also be easily integrated with existing NeRF approaches that rely only on photometric cues as our shadow loss 7 can be used as a regularizer or an auxiliary loss, especially as shadows are invariant to viewpoint changes, surface reflectance properties, or texture changes.

6.2 Conclusions

We show that modern neural rendering techniques can learn neural scene representations (neural shadow fields) and encode 3D geometry just from binary shadow masks. We are motivated by traditional shape-from-X algorithms that typically construct physics-driven inverse models that can exploit cues for 3D reconstruction. We observe that data-driven neural rendering frameworks ignore cues such as shadows, relying on on photometric cues instead. We thus propose a graphics-inspired differentiable shadow rendering component that leverages a volumetric renderer to encode a scene solely from its shadows.

Acknowledgements. This research was supported by the SMART Contract IARPA Grant #2021-20111000004. We would also like to thank Systems & Technology Research (STR). In addition, the authors would also like to thank Professor Voicu Popescu (Purdue University) for being so generous with his time and the valuable discussions that came from our meetings.

References

1. Besl, P., McKay, N.D.: A method for registration of 3-d shapes. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **14**(2), 239–256 (1992). <https://doi.org/10.1109/34.121791>
2. Bobrow, D.G.: Comment on “Numerical shape from shading and occluding boundaries”, pp. 89–94. The MIT Press (1994)
3. Boss, M., Braun, R., Jampani, V., Barron, J.T., Liu, C., Lensch, H.P.: Nerd: Neural reflectance decomposition from image collections. In: *IEEE International Conference on Computer Vision (ICCV)* (2021)
4. Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., Xiao, J., Yi, L., Yu, F.: ShapeNet: An Information-Rich 3D Model Repository. Tech. Rep. arXiv:1512.03012 [cs.GR], Stanford University — Princeton University — Toyota Technological Institute at Chicago (2015)
5. Falcon et al., W.: Pytorch lightning. GitHub. Note: <https://github.com/PyTorchLightning/pytorch-lightning> **3** (2019)
6. Guo, Y., Kang, D., Bao, L., He, Y., Zhang, S.: Nerfren: Neural radiance fields with reflections. *CoRR* **abs/2111.15234** (2021), <https://arxiv.org/abs/2111.15234>
7. Henley, C., Maeda, T., Swedish, T., Raskar, R.: Imaging behind occluders using two-bounce light. In: Vedaldi, A., Bischof, H., Brox, T., Frahm, J.M. (eds.) *Computer Vision – ECCV 2020*. pp. 573–588. Springer International Publishing, Cham (2020)
8. Kato, H., Ushiku, Y., Harada, T.: Neural 3d mesh renderer. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2018)
9. Landabaso, J.L., Pardàs, M., Casas, J.R.: Shape from inconsistent silhouette. *Comput. Vis. Image Underst.* **112**, 210–224 (2008)
10. Li, T.M., Aittala, M., Durand, F., Lehtinen, J.: Differentiable monte carlo ray tracing through edge sampling. *ACM Trans. Graph. (Proc. SIGGRAPH Asia)* **37**(6), 222:1–222:11 (2018)
11. Liu, R., Menon, S., Mao, C., Park, D., Stent, S., Vondrick, C.: Shadows shed light on 3d objects. arXiv e-prints pp. arXiv–2206 (2022)
12. Liu, S., Li, T., Chen, W., Li, H.: Soft rasterizer: A differentiable renderer for image-based 3d reasoning. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 7708–7717 (2019)
13. Lombardi, S., Simon, T., Saragih, J., Schwartz, G., Lehrmann, A., Sheikh, Y.: Neural volumes: Learning dynamic renderable volumes from images. *ACM Trans. Graph.* **38**(4), 65:1–65:14 (Jul 2019)
14. Loper, M.M., Black, M.J.: Opendr: An approximate differentiable renderer. In: *European Conference on Computer Vision*. pp. 154–169. Springer (2014)
15. Lorensen, W.E., Cline, H.E.: Marching cubes: A high resolution 3d surface construction algorithm. *ACM siggraph computer graphics* **21**(4), 163–169 (1987)
16. Martin, W.N., Aggarwal, J.K.: Volumetric descriptions of objects from multiple views. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **PAMI-5**(2), 150–158 (1983). <https://doi.org/10.1109/TPAMI.1983.4767367>
17. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: NeRF: Representing scenes as neural radiance fields for view synthesis. In: *The European Conference on Computer Vision (ECCV)* (2020)
18. Niemeyer, M., Geiger, A.: GIRAFFE: Representing scenes as compositional generative neural feature fields. <https://arxiv.org/abs/2011.12100> (2020)

19. Niemeyer, M., Mescheder, L., Oechsle, M., Geiger, A.: Differentiable volumetric rendering: Learning implicit 3D representations without 3D supervision. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2019)
20. Niemeyer, M., Mescheder, L., Oechsle, M., Geiger, A.: Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (2020)
21. Nimier-David, M., Vicini, D., Zeltner, T., Jakob, W.: Mitsuba 2: A retargetable forward and inverse renderer. *ACM Transactions on Graphics (TOG)* **38**(6), 1–17 (2019)
22. Park, J.J., Florence, P., Straub, J., Newcombe, R., Lovegrove, S.: DeepSDF: Learning continuous signed distance functions for shape representation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 165–174 (2019)
23. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. In: Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., Garnett, R. (eds.) *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc. (2019), <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
24. Quei-An, C.: Nerf_pl: a pytorch-lightning implementation of nerf (2020), https://github.com/kwea123/nerf_pl/
25. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation (2015)
26. Sara Fridovich-Keil and Alex Yu, Tancik, M., Chen, Q., Recht, B., Kanazawa, A.: Plenoxels: Radiance fields without neural networks. In: CVPR (2022)
27. Savarese, S., Rushmeier, H., Bernardini, F., Perona, P.: Shadow carving. In: Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001. vol. 1, pp. 190–197. IEEE (2001)
28. Sitzmann, V., Thies, J., Heide, F., Nießner, M., Wetzstein, G., Zollhofer, M.: Deepvoxels: Learning persistent 3d feature embeddings. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2437–2446 (2019)
29. Srinivasan, P.P., Deng, B., Zhang, X., Tancik, M., Mildenhall, B., Barron, J.T.: Nerv: Neural reflectance and visibility fields for relighting and view synthesis (2020)
30. Tancik, M., Srinivasan, P.P., Mildenhall, B., Fridovich-Keil, S., Raghavan, N., Singhal, U., Ramamoorthi, R., Barron, J.T., Ng, R.: Fourier features let networks learn high frequency functions in low dimensional domains (2020)
31. Tulsiani, S., Efros, A.A., Malik, J.: Multi-view consistency as supervisory signal for learning shape and pose prediction. In: Computer Vision and Pattern Recognition (CVPR) (2018)
32. Velten, A., Willwacher, T., Gupta, O., Veeraraghavan, A., Bawendi, M.G., Raskar, R.: Recovering three-dimensional shape around a corner using ultrafast time-of-flight imaging. *Nature* p. 745 (2012)
33. Verbin, D., Hedman, P., Mildenhall, B., Zickler, T., Barron, J.T., Srinivasan, P.P.: Ref-NeRF: Structured view-dependent appearance for neural radiance fields. arXiv (2021)

34. Vogel, O., Valgaerts, L., Breuß, M., Weickert, J.: Making shape from shading work for real-world images. In: Denzler, J., Notni, G., Süße, H. (eds.) *Pattern Recognition*. pp. 191–200. Springer Berlin Heidelberg, Berlin, Heidelberg (2009)
35. Williams, L.: Casting curved shadows on curved surfaces. In: *Proceedings of the 5th annual conference on Computer graphics and interactive techniques*. pp. 270–274 (1978)
36. Yamazaki, S., Srinivasa Narasimhan, G., Baker, S., Kanade, T.: The theory and practice of coplanar shadowgram imaging for acquiring visual hulls of intricate objects. *International Journal of Computer Vision* **81** (03 2009). <https://doi.org/10.1007/s11263-008-0170-4>
37. Ye, Y., Tulsiani, S., Gupta, A.: Shelf-supervised mesh prediction in the wild. In: *Computer Vision and Pattern Recognition (CVPR)* (2021)
38. Yu, A., Li, R., Tancik, M., Li, H., Ng, R., Kanazawa, A.: PlenOctrees for real-time rendering of neural radiance fields. In: *ICCV* (2021)
39. Zhang, J.Y., Yang, G., Tulsiani, S., Ramanan, D.: NeRS: Neural reflectance surfaces for sparse-view 3d reconstruction in the wild. In: *Conference on Neural Information Processing Systems* (2021)
40. Zhang, R., Tsai, P.S., Cryer, J., Shah, M.: Shape-from-shading: a survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **21**(8), 690–706 (1999). <https://doi.org/10.1109/34.784284>
41. Zheng, Q., Chellappa, R.: Estimation of illuminant direction, albedo, and shape from shading. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **13**(7), 680–702 (1991). <https://doi.org/10.1109/34.85658>