


# Multi-View Photometric Stereo Revisited

**Conference Paper****Author(s):**

Kaya, Berk; Kumar, Suryansh; [Porto de Oliveira, Carlos Eduardo](#) ; Ferrari, Vittorio; Van Gool, Luc

**Publication date:**

2023

**Permanent link:**

<https://doi.org/10.3929/ethz-b-000590876>

**Rights / license:**

[In Copyright - Non-Commercial Use Permitted](#)

**Originally published in:**

<https://doi.org/10.1109/WACV56688.2023.00314>

# Multi-View Photometric Stereo Revisited

Berk Kaya<sup>1</sup> Suryansh Kumar<sup>1\*</sup> Carlos Oliveira<sup>1</sup> Vittorio Ferrari<sup>2</sup> Luc Van Gool<sup>1,3</sup>  
ETH Zürich<sup>1</sup>, Google Research<sup>2</sup>, KU Leuven<sup>3</sup>

## Abstract

*Multi-view photometric stereo (MVPS) is a preferred method for detailed and precise 3D acquisition of an object from images. Although popular methods for MVPS can provide outstanding results, they are often complex to execute and limited to isotropic material objects. To address such limitations, we present a simple, practical approach to MVPS, which works well for isotropic as well as other object material types such as anisotropic and glossy. The proposed approach in this paper exploits the benefit of uncertainty modeling in a deep neural network for a reliable fusion of photometric stereo (PS) and multi-view stereo (MVS) network predictions. Yet, contrary to the recently proposed state-of-the-art, we introduce neural volume rendering methodology for a trustworthy fusion of MVS and PS measurements. The advantage of introducing neural volume rendering is that it helps in the reliable modeling of objects with diverse material types, where existing MVS methods, PS methods, or both may fail. Furthermore, it allows us to work on neural 3D shape representation, which has recently shown outstanding results for many geometric processing tasks. Our suggested new loss function aims to fit the zero level set of the implicit neural function using the most certain MVS and PS network predictions coupled with weighted neural volume rendering cost. The proposed approach shows state-of-the-art results when tested extensively on several benchmark datasets.*

## 1. Introduction

Multi-view photometric stereo (MVPS) aims at recovering accurate and complete 3D reconstruction of an object using multi-view stereo (MVS) and photometric stereo (PS) images [16]. While PS is exemplary in recovering an object’s high-frequency surface details, MVS helps in retaining the global consistency of the object’s 3D shape and assists in correcting overall low-frequency distortion due to PS [10, 23, 34]. Hence, MVPS inherits the complementary output response of PS and MVS methods. Contrary to the active range scanning methods [4, 34, 40], it provides

an efficient, low-cost, and effective alternative for trustworthy 3D data acquisition. And therefore, it is widely preferred in architectural restoration [34], machine vision industry [16, 23, 41], etc.

State-of-the-art geometric methods to solve MVPS indeed provide accurate results but are composed of multiple optimizations and filtering steps applied in sequel [16, 29, 36]. Further, these steps are intricate and require the manual intervention of an expert for precise execution, thereby limiting its automation [29, 36]. Moreover, these approaches cannot meet modern industrial requirements of scalability and low-memory footprint for efficient storage of recovered 3D models. Lately, neural network-based learning methods to solve MVPS have shown few critical advantages over geometric methods [22, 23]. These methods are simpler, effective, and can provide a high-quality 3D model with a lower memory footprint. Yet, they depend on specific assumptions about the material type, which limits their application to anisotropic and glossy material objects.

In this paper, we present a general yet simple and effective approach to the MVPS problem. Inspired by the recent MVPS method [23], we introduce uncertainty modeling in multi-view stereo and photometric stereo neural networks for reliable inference of the 3D position and surface normals, respectively. Although uncertainty estimation helps us filter wrong predictions, it can lead to incomplete recovery of an object’s 3D shape. To this end, Kaya et al. [23] recently proposed Eikonal regularization to recover the missing details due to filtering. On the contrary, we introduce neural volume rendering of the implicit 3D shape representation. It has couple of key advantage over [23] pipeline: (i) It helps extending the application of MVPS to a wider class of object with different material type (see Fig.1(b)). (ii) It further enhances the performance and use of implicit neural shape representation in MVPS leading to state-of-the-art results on benchmark datasets.

Meanwhile, recent multi-view stereo approaches have shown that neural volume rendering using the implicit neural 3D shape representation can effectively model a diverse set of objects via multi-view image rendering techniques [20, 28, 33, 48, 49]. Therefore, introducing it to MVPS can assist in handling challenging objects’ material types. In-

\*Corresponding Author (k.sur46@gmail.com)

tuitively, rendering-based geometry modeling can succeed where both the MVS and PS methods fail to estimate the surface geometry [29, 34, 36]. Further, contrary to the standard practice in MVPS of performing optimization or filtering on explicit geometric primitives [29, 34, 36], *i.e.*, mesh, neural volume rendering relies on neural implicit shape representation, which is memory efficient and is scalable [48]. In summary, our paper makes the following contributions:

- We present a simple, efficient, scalable, and effective MVPS method for the detailed and complete recovery of the object’s 3D shape.
- Our proposed uncertainty-aware neural volume rendering uses confident priors from deep-MVS and deep-PS networks and encapsulates them with an implicit geometric regularizer to solve MVPS demonstrating state-of-the-art reconstruction results on the benchmark dataset [29].
- Contrary to the current state-of-the-art methods, our method applies to a broader class of object material types, including anisotropic and glossy materials. Hence, widen the use of MVPS for 3D data acquisition.

## 2. Related Work

**Classical MVPS.** Early MVPS methods assume a particular analytic BRDF model, which may not be apt for real-world objects whose reflectance differs from the assumed BRDF model [13, 16, 30]. Later, Park *et al.* [36, 37] proposed a piece-wise planar mesh parameterization approach for recovering an object’s fine surface details via displacement texture maps. Nevertheless, their work was not aimed at modeling surface reflectance properties. Other methods such as [8, 39] model the BRDF, yet restricted to near-flat surface modeling assuming the surface normal is known.

Other classical MVPS methods that have been proposed in the last couple of years do provide decent results [29, 51]; yet, their introduced pipeline composes of several complex optimization algorithms such as iso-depth contour estimation, contour tracing, structure-from-motion, multi-view depth propagation, point sorting, mesh optimization using [34], and ACLS algorithm [9]. Moreover, some of these steps require an expert’s intervention for parameter fine-tuning; hence challenging to re-implement, automate and execute. Additionally, the method’s reflectance modeling is built on Alldrin *et al.* [2] and Tan *et al.* [42] work, and therefore, its application is limited to isotropic material objects.

**Deep MVPS.** In recent years, deep learning-based approaches to MVPS have been proposed as alternatives to classical methods. Not long ago, Kaya *et al.* [24] introduced a neural radiance fields-based MVPS approach (NR-MVPS). The proposed pipeline predicts the object’s surface normals using a deep-PS network and blends them in a multi-view volume rendering formulation to solve MVPS. Regardless of its simplicity, it fails to provide a high-quality

3D reconstruction of the object. Further, [22] proposed neural inverse rendering idea to recover an object’s shape and material properties. Among all the deep MVPS methods, the recently introduced uncertainty-based MVPS approach [23] (UA-MVPS) provides better 3D reconstruction results. However, it fails on anisotropic and glossy objects (see Fig.1(b)). On the contrary, this paper proposes a method that can successfully make MVPS 3D acquisition setup work for isotropic, anisotropic, and glossy objects with magnificent results.

## 3. Preliminaries

**MVPS Setup.** Hernández *et al.* [16] proposed the introductory MVPS acquisition setup<sup>1</sup>. It is composed of a turntable arrangement, where light-varying images (PS images) of the object placed on the table are captured from a given viewpoint. Note that the camera and light sources’ position remains fixed, and only the table rotates, providing a new viewpoint ( $v$ ) of the object per rotation. For every table rotation, PS images for each light source are captured and stored (see Fig.1(a)).

**Notation and Definition.** Denoting  $L$  as the total number of point light sources and  $V$  as the total number of viewpoints (corresponds to each table rotation), we define  $\mathcal{X}_{ps}^v = \{X_1^v, X_2^v, \dots, X_L^v\}$  as the set of photometric stereo images from each viewpoint  $v \in [1, V]$ , and  $\mathcal{Y}_{mv} = \{Y^1, Y^2, \dots, Y^V\}$  as the set of multi-view images constructed using  $Y^v = \text{median}(\mathcal{X}_{ps}^v)$  as performed in [29]. The goal of an MVPS algorithm under calibrated setting is to recover the precise and complete geometry of the object. The motivation for using MVS and PS is due to the observation elaborated in [34]. As alluded to above, despite PS can provide reliable high-frequency geometric details, it generally contributes to low-frequency surface distortion at coarse scale [34]. We can correct such distortions using geometric constraints with object’s MVS images.

Using the basic MVPS experimental setup, it is easy to recover two types of surface priors: (i) 3D position per pixel ( $\mathbf{p}_i \in \mathbb{R}^{3 \times 1}$ ) of the object using multi-view stereo images (ii) surface normal for each surface point ( $\mathbf{n}_i^{ps} \in \mathbb{R}^{3 \times 1}$ ) using light varying images [10, 15, 27, 44]<sup>2</sup>. Hence, by design, the problem boils down to effective use MVS and PS surface priors, light varying images, light and camera calibration data for high-quality dense 3D surface recovery. To have the 3D position prior, most methods resort to structure from motion method or its variation [29, 37]. For surface normal prior, one of the popular image formation model is:

<sup>1</sup>Refer Nehab *et al.* [34] 2005 work, which uses active range scanning sensor to solve a similar problem.

<sup>2</sup>Note that MVS reconstruction may not provide reliable per pixel 3D reconstruction. Hence, bad 3D estimates are filtered which leads to sparse set of object 3D points.

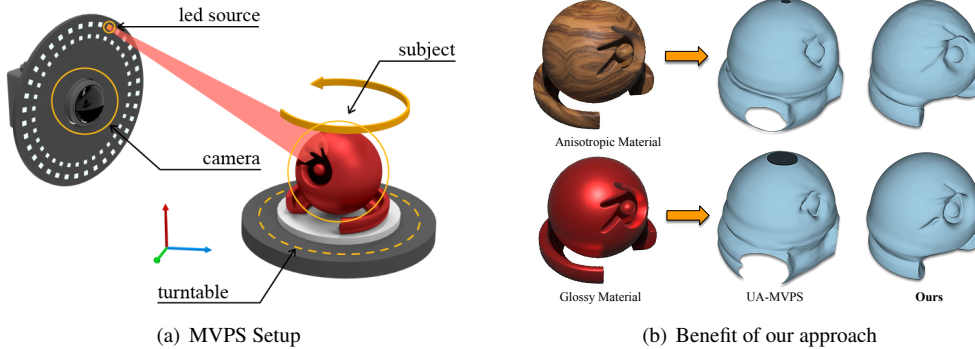


Figure 1. (a) The classical MVPS setup as outlined in Hernández *et al.* [16] work. (b) The advantage of our method over current state-of-the-art deep-MVPS method *i.e.*, UA-MVPS [23]. It can be observed that our method is able to correctly recover the fine object’s details for anisotropic and glossy material object. The 3D model used for the above illustration is taken from [33] dataset.

$$X_j^v(\mathbf{p}_i) = e_j \cdot \rho(\mathbf{n}_i(\mathbf{p}_i), \mathbf{l}_j, \mathbf{v}) \cdot \zeta_a(\mathbf{n}_i(\mathbf{p}_i), \mathbf{l}_j) \cdot \zeta_c(\mathbf{p}_i) \quad (1)$$

Here, the function  $\rho()$  denotes the BRDF,  $\zeta_a(\mathbf{n}_i(\mathbf{p}_i), \mathbf{l}_j) = \max(\mathbf{n}_i(\mathbf{p}_i)^T \mathbf{l}_j, 0)$  accounts for the attached shadow, and  $\zeta_c(\mathbf{p}_i) \in \{0, 1\}$  assigns 0 or 1 value to  $\mathbf{p}_i$  depending on whether it lies in the cast shadow region or not.  $\mathbf{l}_j$  is the light source direction and  $e_j \in \mathbb{R}_+$  is the scalar for light intensity value due to  $j^{th}$  light source. Although surface normal can be estimated with reasonable accuracy using Eq.(1) image formation model [5], modeling BRDF using it can be challenging. Therefore, we propose a neural network-based image rendering approach to overcome such a limitation. Experimental results show that using our approach help MVPS work for a broader class of object material. Next, we describe our approach to the MVPS problem in detail.

## 4. Our Approach

As mentioned in Sec.2, on the one hand, we have the state-of-the-art geometric method that is composed of several complex steps, hence not suitable for automation. Further, it cannot meet the modern demand of scalability, and thus, less convincing for the current challenge of handling a large set of object data. On the other hand, UA-MVPS [23] recent work on deep MVPS is simple and scalable but works well only for isotropic material objects.

This paper proposes a simple, scalable, and effective approach that can handle a much broader range of objects. We first recover the 3D position and surface normal priors from MVS and PS images (MVPS setup) using uncertainty-aware deep multi-view stereo [43] and deep photometric stereo networks [19, 23], respectively. The uncertainty-aware network measures the suitability of the predicted surface measurements for its reliable fusion. However, the filtering of unreliable predictions based on the uncertainty

measures leads to the loss of local surface geometry. Thus, we introduce a geometric regularization term in the overall loss function to recover the complete 3D geometry of the object. To that end, we represent the object’s shape as level sets of a neural network and recover it by optimizing the parameters of a multi-layer perceptron (MLP). The MLP approximates a signed-distance-function (SDF) to a plausible surface based on the point cloud, surface normals, and an implicit geometric regularization term developed on the Eikonal partial differential equation [6].

The above pipeline is inspired by UA-MVPS [23], which generally works well but cannot model anisotropic or glossy surfaces. Hence, not a general solution and is unsuitable for large applications. On a different note, we observed that representing the light fields and density of the object as a neural network in a multi-view volume rendering algorithm improves the 3D reconstruction of general objects. Further, as well-studied, volume rendering generalizes well to diverse objects with different material types. Such an observation leads us to introduce an uncertainty-aware volume rendering approach to the MVPS problem. As we will show, it not only helps achieve state-of-the-art results on isotropic material objects but also provides accurate 3D surface reconstruction on challenging subjects such as glossy texture-less surface objects. Next, we describe each component of our approach in detail, leading to the final loss.

### 4.1. Uncertainty-Aware Deep-MVS Network

Given a set of multi-view images  $\mathcal{Y}_{mv}$ ,  $\{\mathbf{K}_v, \mathbf{R}_v, \mathbf{t}_v\}_{v=1}^V$  the set of camera intrinsics, rotations, and translations for each camera view, the goal is to recover the 3D position of the object corresponding to each pixel with a measure of its reconstruction quality. For that, we use PatchMatch-Net [43] architecture due to its state-of-the-art (SOTA) performance on large-scale images. Further, it provides dense depth maps with per-pixel confidence values. Such an in-

herent property allows the filtering of unreliable depth predictions without having to add an extra uncertainty estimation module into the network.

Built on the idea of classical PatchMatch [3] algorithm, it starts by generating random depth hypotheses. Then, the network repeatedly propagates and evaluates existing depth hypotheses at different image scales in a coarse-to-fine manner. Specifically, feature maps are extracted from each input image and the extracted features are used to generate new depth hypotheses. Subsequently, generated hypotheses are evaluated to compute the matching cost. For that, similarities between warped feature maps are calculated using group-wise correlation [45]. Finally, the depth  $\mathbf{d}_i$  and the confidence  $C_i$  value at pixel  $i$  are computed as follows:

$$\mathbf{d}_i = \sum_{j=1}^{\mathcal{H}} d_i^j \cdot \text{softmax}(\mathbf{J}_i^j), \quad C_i = \text{softmax}(\mathbf{J}_i^{j^*}) \quad (2)$$

Here,  $d_i^j$  is the  $j^{\text{th}}$  depth hypothesis at pixel  $i$  and  $\mathbf{J}_i^j$  is the computed matching cost of corresponding depth hypothesis.  $\mathcal{H}$  is the total number of depth hypotheses, and  $j^*$  is the most likely depth hypothesis at a pixel. After PatchMatchNet is applied at the finest image scale, we obtain the position estimate at pixel coordinates  $\mathbf{o}_i$  by  $\mathbf{p}_i = \mathbf{R}_v(\mathbf{d}_i \mathbf{K}_v^{-1} \mathbf{o}_i) + \mathbf{t}_v$ . Further, we introduce per-pixel binary variable  $c_i^{\text{mvs}}$  to indicate highly confident estimates. We assign  $c_i^{\text{mvs}} = 1$  when  $C_i > \tau_{\text{mvs}}$  and keep  $c_i^{\text{mvs}} = 0$  for the rest [23]. For more details on deep-MVS network’s train and test time specifics refer supplementary or [43].

## 4.2. Uncertainty-Aware Deep-PS Network

To predict surface normals per view from PS images  $\mathcal{X}_{ps}^v$ , and light source directions  $\{\mathbf{l}_j\}_{j=1}^L$ , we use the network architecture presented in [19]. Instead of having a parametric BRDF model assumption, the network learns from training data to map an input *observation map* to a surface normal. An observation map is a 2D matrix-based representation obtained by storing the intensity values at a pixel due to different light sources. Experiments suggest that observation map based representation facilitates accurate estimation of surface normals for general isotropic BRDFs [47, 50]. For more details on the network architecture and observation map refer to Ikehata’s work [19] or supplementary material.

Despite the PS network architecture can predict the object’s surface normals, it cannot measure uncertainty in the predicted value, which is one of the critical components of our approach. Following [23], we adopt the Monte Carlo (MC) dropout approach [11, 12] and build up an uncertainty-aware deep-PS architecture. In a nutshell, we introduce a dropout layer with probability  $p_{\text{mc}}$  after all convolution and fully connected layers. With this adjustment, the network can be treated as a Bayesian neural network, whose parameters approximate a Bernoulli distribu-

tion. Thus, we can train the network with an additional weight decay term scaled by  $\lambda_w$  on network parameters:

$$\mathcal{L}_{ps} = \frac{1}{N_{\text{mc}}} \sum_{j=1}^{N_{\text{mc}}} \|\tilde{\mathbf{n}}_j - \mathbf{n}_{\text{gt}}\|_2^2 + \lambda_w \sum_{k=1}^K \|\mathbf{W}_k\|_2^2 \quad (3)$$

In Eq.(3)  $\tilde{\mathbf{n}}_j$ ,  $\mathbf{n}_{\text{gt}}$  denotes the network’s predicted and ground-truth surface normal, respectively.  $N_{\text{mc}}$  is the number of MC samples and  $\mathbf{W}_k$  stands for the network weights at layer  $k = 1, \dots, K$ . We train the network on CyclesPS dataset [19] once, and used the same network for testing.

At test time, we keep dropout layers active to have a non-deterministic network and we run the network multiple times on the same input. This allows us to capture the fluctuation on the surface normal predictions. We average out all predictions at pixel  $i$  to compute the output normal  $\mathbf{n}_i^{\text{ps}} \in \mathbb{R}^{3 \times 1}$  and the variance  $\tilde{\sigma}_i^2 \in \mathbb{R}^{3 \times 1}$ . Since we are interested in highly confident predictions, we assign  $c_i^{\text{ps}} = 1$  if  $\|\tilde{\sigma}_i^2\|_1 < \tau_{\text{ps}}$  and keep  $c_i^{\text{ps}} = 0$  for the remaining pixels [23]. Here,  $c_i^{\text{ps}}$  is a binary variable to indicate the selection of confident normal prediction.

## 4.3. Shape Representation and Regularization

Using deep-MVS and deep-PS networks —as described above, we filter confident 3D positions and surface normals  $\{\mathbf{p}_i, \mathbf{n}_i^{\text{ps}}\}_{i=1}^{\mathcal{I}} \subset \mathbb{R}^3$  prediction  $\forall i \in [1, \dots, \mathcal{I}]$ . Our goal is to recover object’s dense 3D reconstruction combining those reliable intermediate priors. To this end, we propose to learn the signed distance function (SDF) of the object surface defined by a implicit function  $f_\theta(\mathbf{x}) : \mathbb{R}^3 \rightarrow \mathbb{R}$  using the reliable prediction estimates. We model the function using an MLP parameterized by  $\theta$ , assuming its zero level set approximates the object surface.

To find the optimal  $\theta$ , we consider the Eikonal equation ( $\|\nabla_{\mathbf{x}} f_\theta(\mathbf{x})\| = 1$ ). It establishes a constraint on  $f_\theta(\mathbf{x})$  to represent a true SDF. Note that even if the boundary conditions imposed by the given surface estimates are satisfied (*i.e.*,  $f_\theta(\mathbf{p}_i) = 0$ ,  $\nabla_{\mathbf{x}} f_\theta(\mathbf{p}_i) = \mathbf{n}_i^{\text{ps}}$ ), a unique solution to the zero level set surface may not exist. Nevertheless, describing an incomplete set of surface 3D estimates using Eikonal condition as a regularizer favors smooth and plausible surfaces [14]. Hence, we consider the following regularization term in our optimization:

$$\mathcal{L}_{\text{Eikonal}} = \lambda_e \mathbb{E}_{\mathbf{x}} (\|\nabla_{\mathbf{x}} f_\theta(\mathbf{x})\| - 1)^2 \quad (4)$$

where the expectation is computed w.r.t. a probability distribution  $\mathbf{x} \sim \mathcal{D}$ . Note that recent work [23] has considered the Eikonal regularization to interpolate the surface from MVS and PS network predictions. However, the question we ask in the paper, *did utilize all the imaging prior provided by MVPS well or can we do better?*. In this work, we show that by cleverly using multi-view image prior, we can perform better than UA-MVPS [23]. To accomplish that, we introduce neural volume rendering method to MVPS.

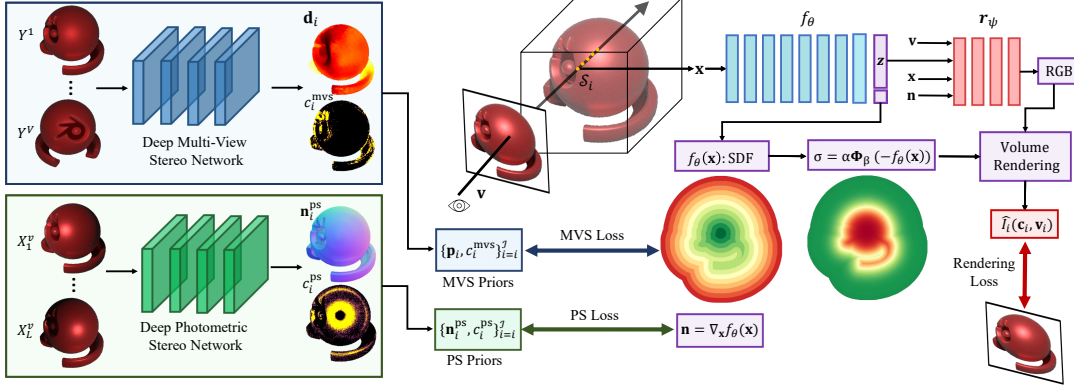


Figure 2. **Method overview (Left to Right)**: We obtain highly confident 3D position and surface normal predictions of the object via uncertainty-aware deep-MVS and deep-PS networks, respectively. Then, we learn the signed distance function representation of the object surface. Finally, our optimization uses the volume rendering technique to recover the missing details of the surface, providing high-quality 3D reconstructions of challenging material types.

#### 4.4. Neural Volume Rendering

Recent work on volume rendering techniques has shown outstanding results in learning scene representations from multi-view images [33]. Although such techniques are impressive with novel view synthesis, they can not faithfully provide the object’s geometry from the learned volume density, leading to inaccurate and noisy reconstructions. Therefore, for our work, we use SDF-based volume rendering approach [48] which models volume density as a function of the signed distance value as follows:

$$\sigma(\mathbf{x}) = \alpha \Phi_\beta(-f_\theta(\mathbf{x})),$$

$$\text{where } \Phi_\beta(s) = \begin{cases} \frac{1}{2} \exp(\frac{s}{\beta}), & \text{if } s \leq 0 \\ 1 - \frac{1}{2} \exp(-\frac{s}{\beta}), & \text{if } s > 0 \end{cases} \quad (5)$$

Here,  $\alpha, \beta > 0$  are trainable parameters and  $\Phi_\beta(\cdot)$  is the cumulative distribution function of a zero-mean Laplace distribution. Eq:(5) ensures a smooth transition of density values near the object boundary, and at the same time allows a suitable extraction of zero level set after optimization for surface recovery. Inspired by the classical volume rendering techniques [21, 32], the expected color  $I(\mathbf{c}_i, \mathbf{v}_i)$  of a camera ray  $\mathbf{x}_i(t) = \mathbf{c}_i + t\mathbf{v}_i$  with camera center  $\mathbf{c}_i \in \mathbb{R}^3$  and viewing direction vector  $\mathbf{v}_i \in \mathbb{R}^3$  can be modeled as:

$$I(\mathbf{c}_i, \mathbf{v}_i) = \int_{t_n}^{t_f} T(\mathbf{x}_i(t)) \sigma(\mathbf{x}_i(t)) \mathbf{r}_\psi(\mathbf{x}_i(t), \mathbf{n}_i(t), \mathbf{v}_i) dt, \quad (6)$$

where  $T(\mathbf{x}_i(t)) = \exp(-\int_0^t \sigma(\mathbf{x}_i(s)) ds)$  is the transparency,  $\mathbf{n}_i(t) = \nabla_{\mathbf{x}} f_\theta(\mathbf{x}_i(t))$  is the level set’s normal at  $\mathbf{x}_i(t)$ ,  $\mathbf{r}_\psi$  is the radiance field function and  $(t_n, t_f)$  are the bounds of the ray. Using the quadrature rule for numerical integration [32] and the ray sampling strategy in [49], we

approximate the expected color as :

$$\hat{I}(\mathbf{c}_i, \mathbf{v}_i) = \sum_{j \in \mathcal{S}_i} T_j (1 - \exp(-\sigma_j \delta_j)) \mathbf{r}_\psi(\mathbf{x}_j, \mathbf{n}_j, \mathbf{v}) \quad (7)$$

Here,  $\mathcal{S}_i$  is the set of samples along the ray,  $\delta_j$  is the distance between each adjacent samples and  $T_j$  is the approximated transparency [49]. To realize  $\mathbf{r}_\psi$ , we introduce a second MLP with learnable parameters  $\psi$ . The radiance fields network  $\mathbf{r}_\psi$  is placed subsequent to the signed distance field network  $f_\theta$  (see Fig. 2). Furthermore, we introduce a feature vector  $\mathbf{z} \in \mathbb{R}^{256}$  that is extracted from  $f_\theta$  using a fully connected layer. This feature vector is fed to the radiance field network  $\mathbf{r}_\psi$  to account for global illumination effects. We optimize  $f_\theta$  and  $\mathbf{r}_\psi$  network on the test subject together. After optimization, we extract the zero level set of  $f_\theta$  and recover the shape mesh using marching cubes algorithm [31]. For more details, refer to Sec.§5.1 and [48].

**Optimization.** Our overall training loss is as follows:

$$\mathcal{L}_{\text{mvps}} = \frac{1}{\mathcal{I}} \sum_{i=1}^{\mathcal{I}} \left( \overbrace{c_i^{\text{mvs}} |f_\theta(\mathbf{p}_i)|}^{\text{MVS Loss}} + \overbrace{c_i^{\text{ps}} \|\mathbf{n}_i^r - \mathbf{n}_i^{\text{ps}}\|}^{\text{PS Loss}} \right) + \underbrace{(1 - c_i^{\text{mvs}} c_i^{\text{ps}}) \|I_i - \hat{I}(\mathbf{c}_i, \mathbf{v}_i)\|_1}_{\text{Rendering Loss}} + \underbrace{\lambda_m \sum_{i \in \mathcal{M}} CE(\max(\sigma_j / \alpha), 0)}_{\text{Mask Loss}} + \underbrace{\lambda_e \mathbb{E}_{\mathbf{x}} (\|\nabla_{\mathbf{x}} f_\theta(\mathbf{x})\| - 1)^2}_{\text{Eikonal Regularization}} \quad (8)$$

Eq.(8) consists of five terms. Here, the first term forces the signed distance to vanish on the high fidelity position predictions of deep-MVS network. Similarly, the second term encourages the expected surface normal on a ray  $\mathbf{n}_i^r = \sum_{j \in \mathcal{S}_i} T_j (1 - \exp(-\sigma_j \delta_j)) \mathbf{n}_i(t)$  to align with the

highly confident deep-PS predictions. The third term introduces an uncertainty-aware rendering loss to the optimization for the pixels where either MVS or PS fails. Intuitively, this allows the optimization to recover the missing surface details using rendering. We further improve the geometry using the object masks. For that, we first find the maximum density on rays outside the object mask (i.e.  $i \in \mathcal{M}$ ). Then, we apply cross-entropy loss (CE) to minimize ray and geometry intersections as in [49]. The final term applies Eikonal regularization for plausible surface recovery as discussed in Sec. §4.3. Fig. (2) shows the overall pipeline of our proposed approach.

## 5. Experiment and Results

**Datasets.** First, we evaluated our approach on the DiLiGenT-MV [29]. DiLiGenT-MV is a standard benchmark for the MVPS setup, consisting of five real-world objects. The images are acquired using a turntable setup where the object is placed  $\sim 1.5m$  away from the camera. The turntable is rotated with 20 uniform rotations for each object, and 96 distinct light sources are used to capture light-varying images at each rotation. Although the DiLiGenT-MV benchmark consists of challenging objects with non-Lambertian surfaces, all provided objects satisfy isotropic BRDF property. Therefore, we simulated a new dataset consisting for objects with anisotropic and glossy surfaces.

Similar to classical setup, we simulated our dataset using a turntable setup with 36 angle rotations. We place 72 light sources in a concentric way around the camera (see Fig.1(a)) and rendered images corresponding to each light source. We use licensed Houdini software to simulate our setup and render MVPS images of a single object 3D model taken from NeRF synthetic dataset [33] with three different material types (Wood, Gray, Red)<sup>3</sup>. The Wood category is rendered to study anisotropic material behavior and the other two categories to analyse our method’s performance on texture-less glossy objects. We rendered images at  $1280 \times 720$  resolution to better capture the object details<sup>4</sup>.

### 5.1. Implementation Details

We implemented our method in Python 3.8 using PyTorch 1.7.1 [38] and conducted all our experiments on a single NVIDIA GPU with 11GB of RAM. We first train uncertainty-aware deep-MVS and deep-PS networks under a supervised setting. Then, we use these networks to have 3D position and surface normal predictions at test time. Finally, MVS images, along with the network predictions and their per-pixel confidence values, are used to optimize the proposed loss function (Eq.(8)).

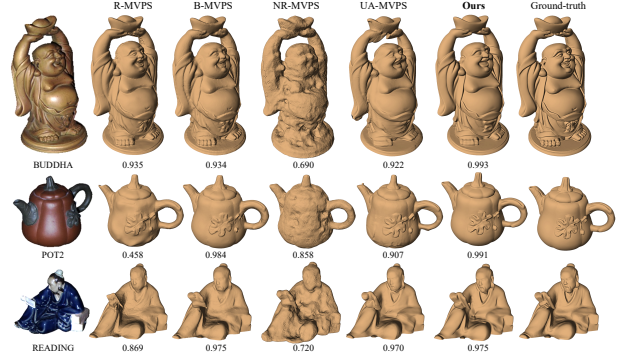


Figure 3. Comparison of MVPS reconstructions on DiLiGenT-MV benchmark [29]. We report  $F$ -score metric results for numerical comparison. We can observe that our method recovers fine details and provides high-quality reconstructions of challenging objects.

**(a) Deep-MVS Network.** The deep-MVS network is trained on DTU’s train set [1]. The training takes 8 epochs using the learning rate 0.001 and Adam optimizer [25]. We use the MVS trained model at three coarser stages at test time to predict depth with coarse-to-fine approach. The depth  $d_i$  and the confidence  $C_i$  at each pixel  $i$  are computed using Eq:(2). The predicted depth is further enhanced using [18] work and converted to a set of 3D points  $\{\mathbf{p}_i\}_{i=1}^T$  by back-projecting the depth values to 3D space. Finally, we obtain binary confidences  $c_i^{mvs}$  by setting  $\tau_{mvs} = 0.9$  for reliable fusion of confident position predictions.

**(b) Deep-PS Network.** We train the deep-PS network on CyclesPS dataset [19] for 10 epochs using Adam optimizer [25] and learning rate of 0.1. We use probability of  $p_{mc} = 0.2$  in every dropout layer of the architecture. For training, we set  $N_{mc} = 10$  and  $\lambda_w = 10^{-4}$  (see Eq:(3)). At test time, we first create observation map per-pixel using MVPS images. We then run the network on each observation map 100 times to have the output surface normal  $\mathbf{n}_i^{ps}$  and the prediction variance  $\hat{\sigma}_i^2$  [11, 12]. Finally, we obtain the confidence value  $c_i^{ps}$  at  $i^{th}$  pixel by setting  $\tau_{ps} = 0.03$ .

**(c) Overall Shape Optimization.** As described in §Sec.4.4, we optimize two networks during optimization: signed distance field network ( $f_\theta$ ) and radiance field network ( $\mathbf{r}_\psi$ ).  $f_\theta$  consists of 8 MLP layers with a skip connection connecting the first layer to the 4<sup>th</sup>. On the other hand,  $\mathbf{r}_\psi$  has four MLP layers (see Fig.2). All the layers of both networks have 256 units. We apply Fourier feature encoding to the inputs (position  $\mathbf{x}$  and view direction  $\mathbf{v}$ ) to improve the networks’ ability to represent high-frequency details [33]. For the loss function in Eq:(8), we set  $\lambda_m = 0.1$  and  $\lambda_e = 1$ . We use a set of multi-view images which are captured under the illumination of the same randomly chosen light source to compute the rendering loss. We use Adam optimizer [25] with learning rate  $10^{-4}$  and train for  $10^4$  epochs. In each epoch, we use batches of 1024 rays

<sup>3</sup>CC-BY-3.0 license.

<sup>4</sup>Our dataset and further details related to it will be available soon.

Method Category →		Deep Multi-View Stereo		Photometric Stereo			View-Synthesis		Ours
Dataset ↓	Method →	MVSNet [46]	PM-Net [43]	Robust PS [35]	SDPS-Net [5]	CNN-PS [19]	NeRF [33]	VolSDF [48]	
	BEAR	0.135	0.672	0.266	0.239	0.293	0.865	0.962	<b>0.965</b>
	BUDDHA	0.147	0.799	0.367	0.298	0.363	0.713	0.786	<b>0.993</b>
	COW	0.095	0.734	0.245	0.447	0.511	0.810	0.985	<b>0.987</b>
	POT2	0.126	0.666	0.231	0.464	0.632	0.859	0.946	<b>0.991</b>
	READING	0.115	0.834	0.242	0.188	0.508	0.673	0.683	<b>0.975</b>
	<b>AVERAGE</b>	0.124	0.741	0.270	0.327	0.461	0.784	0.873	<b>0.982</b>

Table 1.  $F$ -score comparison of standalone method reconstructions on DiLiGenT-MV benchmark [29]. Our method outperforms standalone multi-view stereo, photometric stereo and view synthesis methods in all of the object categories.

from each view and sample 64 points along each ray [48]. To compute the Eikonal regularization as in Eq.(4), we also uniformly sample points globally. So, the distribution  $\mathcal{D}$  stands for the collection of these ray samples and global samples. After the optimization, we extract zero level set of the learned SDF representation by  $f_\theta$  and recover the shape mesh using marching cubes algorithm [31] on a  $512^3$  grid.

## 5.2. Statistical Analysis

We performed comparative analysis on the DiLiGenT-MV dataset [29]. To evaluate the quality of the shape reconstructions, we use well-known Chamfer- $L_2$  and  $F$ -score [26] metric. For better understanding, we present the performance comparison result in two different categories depending on the method type.

**(a) Standalone Method Comparison.** By the standalone method, we refer to the approaches that use only one modality *i.e.*, either MVS or PS images for 3D reconstruction. We consider SOTA MVS, PS, and view-synthesis methods for this comparison. Note that we use Horn and Brooks algorithm [17] for normal integration to recover depth maps. We then back-project the recovered depths to 3D space to evaluate reconstruction performance. Table 1 presents the  $F$ -score comparison of these methods on DiLiGenT-MV [29]. The statistics show that our method consistently outperforms the standalone approaches. Further, we observed that none of the standalone methods could reliably recover the object’s 3D shape. On the contrary, our method gives accurate reconstruction by effectively exploiting the complementary surface and image priors.

**(b) MVPS Methods Comparison.** Table 2 provides the  $F$ -score comparison results with SOTA MVPS methods on the DiLiGenT-MV benchmark dataset. For our comparison, we consider both explicit geometry modeling-based classical approaches [29, 36], and neural implicit representation based deep approaches [23, 24]. The numerical results show that our method provides the highest scores on three objects categories. Moreover, it outperforms all the existing MVPS methods on average. Some important point to note is that **(i)** Our approach provides a scalable and easy-to-execute implementation, without requiring tedious sequential steps as in classical methods [29], **(ii)** Our MLP based shape rep-

Dataset ↓	Method →	R-MVPS [36]	B-MVPS [29]	NR-MVPS [24]	UA-MVPS [23]	Ours
	BEAR	0.504	<b>0.986</b>	0.856	0.895	0.965
	BUDDHA	0.935	0.934	0.690	0.922	<b>0.993</b>
	COW	0.915	<b>0.989</b>	0.844	0.979	0.987
	POT2	0.458	0.984	0.858	0.907	<b>0.991</b>
	READING	0.869	0.975	0.720	0.970	<b>0.975</b>
	<b>AVERAGE</b>	0.736	0.974	0.794	0.935	<b>0.982</b>

Table 2.  $F$ -score comparison of MVPS reconstructions on DiLiGenT-MV benchmark [29]. Our method performs consistently well on various objects and is better than others on average.

resentation requires only 3.07MB of memory, while explicit geometric methods may require up to 90MB. Such advantages make our method an efficient and effective algorithmic choice for solving MVPS.

## 5.3. Further Analysis

**(a) Anisotropic and Textureless Glossy Surfaces.** We perform evaluations on our synthetic dataset to analyze the efficiency of our approach on anisotropic and texture-less glossy surfaces. In Fig.4(a), we provide Chamfer  $L_2$  metric comparison of our method with the recent UA-MVPS [23]. The results show that our method performs much better than its competitor on glossy (Gray, Red) and anisotropic surfaces (Wood). In Fig.4(b), we show qualitative results of the uncertainty-aware deep-MVS and deep-PS networks on the Gray category. It can be observed from visual results that deep-MVS cannot provide reliable position estimates on texture-less glossy surfaces. For this reason, methods relying on the fusion of only MVS and PS priors (such as UA-MVPS) cannot handle all kinds of surfaces. On the other hand, our method can recover the missing surface information by effectively utilizing volume rendering; hence, it can suitably work for anisotropic and glossy surface profiles.

**(b) Optimization.** Here, we investigate the effectiveness of our proposed optimization loss in Eq.(8) with an ablation study. For that, we compare the reconstruction quality of our method by removing *(i)* MVS loss term, *(ii)* PS loss term, *(iii)* rendering loss term and *(iv)* uncertainty modeling ( $c_i^{mvs}$  and  $c_i^{ps}$ ) from the overall loss. In Table 3, we provide Chamfer  $L_2$  metric comparison of the reconstruction quality achieved under each of these configurations. The numerical results verify that uncertainty modeling based integration of MVS, PS and rendering loss terms provides best



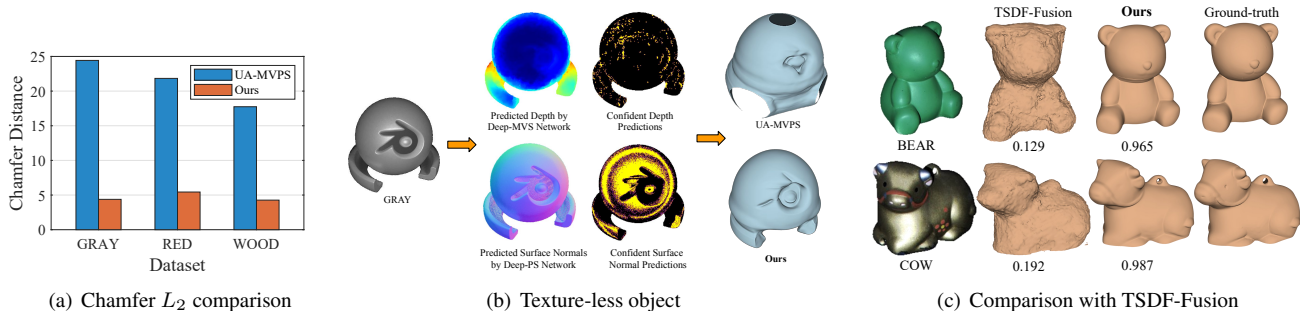


Figure 4. (a) Chamfer  $L_2$  comparison of our method with UA-MVPS [23] on our synthetic dataset (lower is better). (b) We show depth and surface normal predictions on texture-less object. Pixels marked with yellow color indicate confident MVS or PS predictions ( $c_i^{mvs}$  and  $c_i^{ps}$ ). Note that MVS cannot predict depth reliably on texture-less surface, which leads to inferior results in UA-MVPS [23]. On the other hand, our uncertainty-aware volume rendering approach can recover missing surface information, and therefore, provides better reconstructions. (c) Comparison of our method with TSDF Fusion algorithm [7]. We report  $F$ -score metric for numerical comparison.

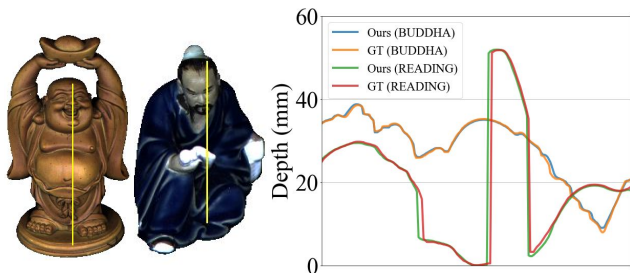


Figure 5. Surface profile of our reconstructions on a randomly chosen path. Clearly, our surface profile overlaps with the ground truth (GT), which indicates the high quality of our reconstructions.

results on DiLiGenT-MV [29].

(c) **Surface Profile.** To show the quality of our recovered 3D reconstructions, we study the surface topology across an arbitrarily chosen curve on the surface. Fig.5 shows a couple of examples of such surface profile on Buddha and Cow sequences. Clearly, our recovered surface profiles align well with the ground truth.

(d) **Volumetric Fusion Approach.** Of course, one can use robust 3D fusion method such as TSDF fusion [7] to recover the object’s 3D reconstruction. And therefore, we conducted this experiment to study the results that can be recovered using such fusion techniques. Accordingly, we fuse deep-MVS depth and the depth from deep-PS normal integration [17] using the TSDF fusion. Fig.4(c) shows that TSDF fusion provide inferior results compared to ours.

(e) **Limitations.** Although our method works well on glossy objects, it may fail on materials with mirror reflection. Furthermore, SDF representation of the object shape restricts our approach to solid and opaque materials. Finally, our work considers a calibrated setting for MVPS setup, and it would be interesting to further investigate our

Settings   Dataset →	BEAR	BUDDHA	COW	POT2	READING	AVERAGE
w/o MVS Loss	0.189	0.089	0.202	0.156	0.353	0.198
w/o PS Loss	0.301	0.572	0.184	0.262	0.428	0.349
w/o Rendering Loss	0.154	0.471	0.269	0.235	0.374	0.301
w/o Uncertainty-Aware.	0.267	0.085	0.313	0.137	0.251	0.211
<b>Ours</b>	<b>0.213</b>	<b>0.088</b>	<b>0.176</b>	<b>0.198</b>	<b>0.253</b>	<b>0.186</b>

Table 3. Contribution of MVS, PS, rendering loss terms and uncertainty modeling to our reconstruction quality. We report Chamfer  $L_2$  metric for comparison (lower is better). Clearly, our proposed loss in Eq.(8) produces best results on average.

approach in an uncalibrated setup. *For more results and exhaustive analysis of our method refer to our supplementary.*

## 6. Conclusion

The proposed method addresses the current limitations of well-known MVPS methods and makes it work well for diverse object material types. Experimental studies on anisotropic and texture-less glossy objects show that existing MVS and PS modeling techniques may not always extract essential cues for accurate 3D reconstructions. However, by integrating incomplete yet reliable MVS and PS information into a rendering pipeline and leveraging the generalization ability of the modern view synthesis approach to model complex BRDFs, it is possible to make MVPS setup work well for anisotropic materials and glossy texture-less objects with better accuracy. Finally, the performance on the standard benchmark shows that our method outperforms existing methods providing exemplary 3D reconstruction results. To conclude, we believe that our approach will open up new avenues for applying MVPS to real-world applications such as metrology, forensics, etc.

**Acknowledgement.** The authors thank ETH support to the Computer Vision Lab (CVL) and Focused Research Award from Google (ETH 2019-HE-318, 2019-HE-323, 2020-FS-351, 2020-HS-411).

## References

- [1] Henrik Aanæs, Rasmus Ramsbøl Jensen, George Vogiatzis, Engin Tola, and Anders Bjorholm Dahl. Large-scale data for multiple-view stereopsis. *International Journal of Computer Vision*, 120(2):153–168, 2016. [6](#)
- [2] Neil Alldrin, Todd Zickler, and David Kriegman. Photometric stereo with non-parametric and spatially-varying reflectance. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2008. [2](#)
- [3] Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. Patchmatch: A randomized correspondence algorithm for structural image editing. *ACM Trans. Graph.*, 28(3):24, 2009. [4](#)
- [4] Avishek Chatterjee and Venu Madhav Govindu. Efficient and robust large-scale rotation averaging. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 521–528, 2013. [1](#)
- [5] Guanying Chen, Kai Han, Boxin Shi, Yasuyuki Matsushita, and Kwan-Yee K Wong. Self-calibrating deep photometric stereo networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8739–8747, 2019. [3](#), [7](#)
- [6] Michael G Crandall and Pierre-Louis Lions. Viscosity solutions of hamilton-jacobi equations. *Transactions of the American mathematical society*, 277(1):1–42, 1983. [3](#)
- [7] Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 303–312, 1996. [8](#)
- [8] Yue Dong, Jiaping Wang, Xin Tong, John Snyder, Yanxiang Lan, Moshe Ben-Ezra, and Baining Guo. Manifold bootstrapping for svbrdf capture. *ACM Transactions on Graphics (TOG)*, 29(4):1–10, 2010. [2](#)
- [9] Lawrence et al. Inverse shade trees for non-parametric material representation and editing. *ACM Transactions on Graphics (TOG)*, pages 735–745, 2006. [2](#)
- [10] Yasutaka Furukawa and Carlos Hernández. Multi-view stereo: A tutorial. *Foundations and Trends® in Computer Graphics and Vision*, 9(1-2):1–148, 2015. [1](#), [2](#)
- [11] Yarín Gal and Zoubin Ghahramani. Bayesian convolutional neural networks with bernoulli approximate variational inference. *arXiv preprint arXiv:1506.02158*, 2015. [4](#), [6](#)
- [12] Yarín Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016. [4](#), [6](#)
- [13] Dan B Goldman, Brian Curless, Aaron Hertzmann, and Steven M Seitz. Shape and spatially-varying brdfs from photometric stereo. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(6):1060–1071, 2009. [2](#)
- [14] 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. [4](#)
- [15] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003. [2](#)
- [16] Carlos Hernandez, George Vogiatzis, and Roberto Cipolla. Multiview photometric stereo. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(3):548–554, 2008. [1](#), [2](#), [3](#)
- [17] Berthold KP Horn and Michael J Brooks. The variational approach to shape from shading. *Computer Vision, Graphics, and Image Processing*, 33(2):174–208, 1986. [7](#), [8](#)
- [18] Tak-Wai Hui, Chen Change Loy, and Xiaoou Tang. Depth map super-resolution by deep multi-scale guidance. In *European conference on computer vision*, pages 353–369. Springer, 2016. [6](#)
- [19] Satoshi Ikehata. Cnn-ps: Cnn-based photometric stereo for general non-convex surfaces. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–18, 2018. [3](#), [4](#), [6](#), [7](#)
- [20] Nishant Jain, Suryansh Kumar, and Luc Van Gool. Robustifying the multi-scale representation of neural radiance fields. *arXiv preprint arXiv:2210.04233*, 2022. [1](#)
- [21] James T Kajiya and Brian P Von Herzen. Ray tracing volume densities. *ACM SIGGRAPH computer graphics*, 18(3):165–174, 1984. [5](#)
- [22] Berk Kaya, Suryansh Kumar, Carlos Oliveira, Vittorio Ferrari, and Luc Van Gool. Uncalibrated neural inverse rendering for photometric stereo of general surfaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3804–3814, 2021. [1](#), [2](#)
- [23] Berk Kaya, Suryansh Kumar, Carlos Oliveira, Vittorio Ferrari, and Luc Van Gool. Uncertainty-aware deep multi-view photometric stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12601–12611, 2022. [1](#), [2](#), [3](#), [4](#), [7](#), [8](#)
- [24] Berk Kaya, Suryansh Kumar, Francesco Sarno, Vittorio Ferrari, and Luc Van Gool. Neural radiance fields approach to deep multi-view photometric stereo. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1965–1977, 2022. [2](#), [7](#)
- [25] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. [6](#)
- [26] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG)*, 36(4):1–13, 2017. [7](#)
- [27] Kiriakos N Kutulakos and Steven M Seitz. A theory of shape by space carving. *International journal of computer vision*, 38(3):199–218, 2000. [2](#)
- [28] Soomin Lee, Le Chen, Jiahao Wang, Alexander Liniger, Suryansh Kumar, and Fisher Yu. Uncertainty guided policy for active robotic 3d reconstruction using neural radiance fields. *IEEE Robotics and Automation Letters*, 2022. [1](#)
- [29] Min Li, Zhenglong Zhou, Zhe Wu, Boxin Shi, Changyu Diao, and Ping Tan. Multi-view photometric stereo: A robust solution and benchmark dataset for spatially varying isotropic materials. *IEEE Transactions on Image Processing*, 29:4159–4173, 2020. [1](#), [2](#), [6](#), [7](#), [8](#)

- [30] Jongwoo Lim, Jeffrey Ho, Ming-Hsuan Yang, and David Kriegman. Passive photometric stereo from motion. In *Tenth IEEE International Conference on Computer Vision (ICCV'05) Volume 1*, volume 2, pages 1635–1642. IEEE, 2005. [2](#)
- [31] 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. [5](#), [7](#)
- [32] Nelson Max. Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics*, 1(2):99–108, 1995. [5](#)
- [33] 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 *European conference on computer vision*, pages 405–421. Springer, 2020. [1](#), [3](#), [5](#), [6](#), [7](#)
- [34] Diego Nehab, Szymon Rusinkiewicz, James Davis, and Ravi Ramamoorthi. Efficiently combining positions and normals for precise 3D geometry. *ACM Transactions on Graphics (Proc. of ACM SIGGRAPH 2005)*, 24(3):536–543, 2005. [1](#), [2](#)
- [35] Tae-Hyun Oh, Hyeongwoo Kim, Yu-Wing Tai, Jean-Charles Bazin, and In So Kweon. Partial sum minimization of singular values in rpca for low-level vision. In *Proceedings of the IEEE international conference on computer vision*, pages 145–152, 2013. [7](#)
- [36] Jaesik Park, Sudipta N Sinha, Yasuyuki Matsushita, Yu-Wing Tai, and In So Kweon. Robust multiview photometric stereo using planar mesh parameterization. *IEEE transactions on pattern analysis and machine intelligence*, 39(8):1591–1604, 2016. [1](#), [2](#), [7](#)
- [37] Jaesik Park, Sudipta N Sinha, Yasuyuki Matsushita, Yu-Wing Tai, and In So Kweon. Multiview photometric stereo using planar mesh parameterization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1161–1168, 2013. [2](#)
- [38] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. [6](#)
- [39] Peiran Ren, Jiaping Wang, John Snyder, Xin Tong, and Baining Guo. Pocket reflectometry. *ACM Transactions on Graphics (TOG)*, 30(4):1–10, 2011. [2](#)
- [40] Erik Sandström, Martin R Oswald, Suryansh Kumar, Silvan Weder, Fisher Yu, Cristian Sminchisescu, and Luc Van Gool. Learning online multi-sensor depth fusion. *arXiv preprint arXiv:2204.03353*, 2022. [1](#)
- [41] Francesco Sarno, Suryansh Kumar, Berk Kaya, Zhiwu Huang, Vittorio Ferrari, and Luc Van Gool. Neural architecture search for efficient uncalibrated deep photometric stereo. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 361–371, 2022. [1](#)
- [42] Ping Tan, Satya P Mallick, Long Quan, David J Kriegman, and Todd Zickler. Isotropy, reciprocity and the generalized bas-relief ambiguity. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2007. [2](#)
- [43] Fangjinhua Wang, Silvano Galliani, Christoph Vogel, Pablo Speciale, and Marc Pollefeys. Patchmatchnet: Learned multi-view patchmatch stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14194–14203, 2021. [3](#), [4](#), [7](#)
- [44] Robert J Woodham. Photometric method for determining surface orientation from multiple images. *Optical engineering*, 19(1):191139, 1980. [2](#)
- [45] Qingshan Xu and Wenbing Tao. Learning inverse depth regression for multi-view stereo with correlation cost volume. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 12508–12515, 2020. [4](#)
- [46] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 767–783, 2018. [7](#)
- [47] Zhuokun Yao, Kun Li, Ying Fu, Haofeng Hu, and Boxin Shi. Gps-net: Graph-based photometric stereo network. *Advances in Neural Information Processing Systems*, 33, 2020. [4](#)
- [48] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems*, 34, 2021. [1](#), [2](#), [5](#), [7](#)
- [49] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems*, 33, 2020. [1](#), [5](#), [6](#)
- [50] Qian Zheng, Yiming Jia, Boxin Shi, Xudong Jiang, Ling-Yu Duan, and Alex C Kot. Spline-net: Sparse photometric stereo through lighting interpolation and normal estimation networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 8549–8558, 2019. [4](#)
- [51] Zhenglong Zhou, Zhe Wu, and Ping Tan. Multi-view photometric stereo with spatially varying isotropic materials. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1482–1489, 2013. [2](#)