

PANDORA: Polarization-Aided Neural Decomposition Of Radiance

Akshat Dave, Yongyi Zhao and Ashok Veeraraghavan
ECE Department
Rice University, Houston, USA

Abstract

Reconstructing an object’s geometry and appearance from multiple images, also known as inverse rendering, is a fundamental problem in computer graphics and vision. Inverse rendering is inherently ill-posed because the captured image is an intricate function of unknown lighting conditions, material properties and scene geometry. Recent progress in representing scene properties as coordinate-based neural networks have facilitated neural inverse rendering resulting in impressive geometry reconstruction and novel-view synthesis. Our key insight is that polarization is a useful cue for neural inverse rendering as polarization strongly depends on surface normals and is distinct for diffuse and specular reflectance. With the advent of commodity, on-chip, polarization sensors, capturing polarization has become practical. Thus, we propose PANDORA, a polarimetric inverse rendering approach based on implicit neural representations. From multi-view polarization images of an object, PANDORA jointly extracts the object’s 3D geometry, separates the outgoing radiance into diffuse and specular and estimates the illumination incident on the object. We show that PANDORA outperforms state-of-the-art radiance decomposition techniques. PANDORA outputs clean surface reconstructions free from texture artefacts, models strong specularities accurately and estimates illumination under practical unstructured scenarios.

Keywords— Polarization, inverse rendering, multi-view reconstruction, implicit neural representations

1 Introduction

Inverse rendering involves reconstructing an object’s appearance and geometry from multiple images of the object captured under different viewpoints and/or lighting conditions. It is important for many computer graphics and vision applications such as re-lighting, synthesising novel views and blending real objects with virtual scenes. Inverse rendering is inherently challenging because the object’s 3D shape, surface reflectance and incident illumination are intermixed in the captured images. A diverse array of techniques have been proposed to alleviate this challenge by incorporating prior knowledge about the scene, by optimizing the scene parameters iteratively using differentiable rendering and by using imaging modalities that exploit unique properties of light such as spectrum, polarization and time.

Neural Inverse Rendering. Recent works demonstrate that modelling the outgoing radiance and object shape as coordinate-based neural networks results in impressive novel-view synthesis (NeRF) [28] and surface reconstruction (VolSDF) [52] from multi-view captures. The outgoing radiance from the object is a combination of different components of surface reflectance and illumination incident on the object. As a result, separation and modification of components of the captured object’s reflectance is not possible with works such as NeRF and VolSDF. Moreover, the diffuse and specular components of object reflectance have different view dependence. Using the same network to model a combination of diffuse and specular radiance results in inaccurate novel view synthesis.

*Email: akshat@rice.edu, Project Webpage: akshatdave.github.io/pandora

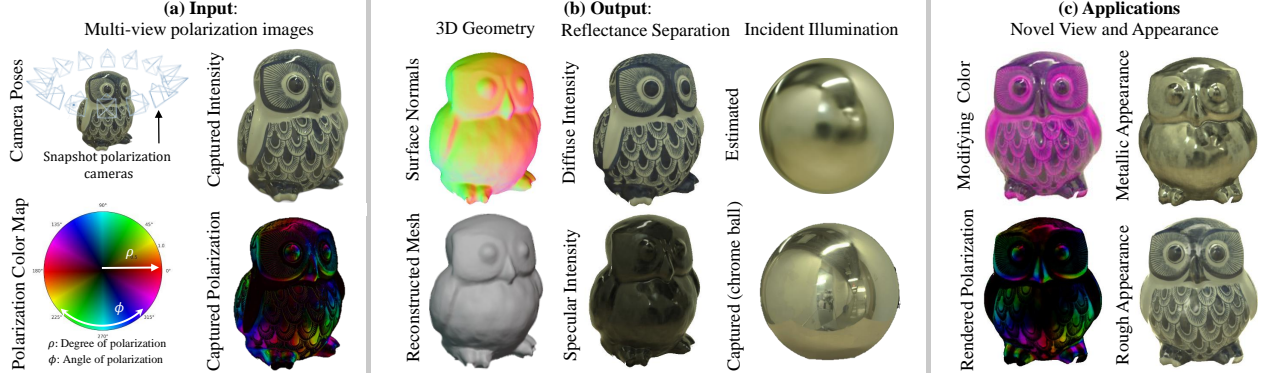


Figure 1: **PANDORA Overview:** PANDORA utilizes multi-view polarization images with known poses (a) and outputs the object’s 3D geometry, separation of radiance in to diffuse and specular along with incident illumination (b). The learned PANDORA model can be applied to render the object under novel views and edit the object’s appearance (c). Please refer to our project webpage for renderings of these outputs and applications under varying viewpoints.

Radiance Decomposition. Decomposition of the outgoing radiance into reflectance parameters and incident illumination is inherently ill-posed. Recent works such as PhySG [56] and NeuralPIL [8] aim to address the ill-posed nature of radiance decomposition by employing spherical Gaussians and data-driven embeddings respectively to model the reflectance and lighting. While these approach provides plausible decomposition in simple settings under large number of measurements, the decomposition is inaccurate in challenging scenarios such as strong specularities and limited views (Fig. 4) leading to blurrier specular reconstructions and artefacts in surface reconstruction.

Key idea: Polarization as a cue for reflectance decomposition. Our key insight is that polarimetric cues aid in radiance decomposition. Polarization has a strong dependence on surface normals. The diffuse and specular reflectance components have different polarimetric properties: the specular is more polarized than diffuse and the polarization angle for the two components are orthogonal. The advent of snapshot polarimetric cameras has made it pratical to capture this polarization information. In this work, we present our approach PANDORA that exploits multi-view polarization images for jointly recovering the 3D surface, separating the diffuse-specular radiance and estimating the incident illumination.

Our approach. PANDORA models the geometry as an implicit neural surface similar to VolSDF. Implicit coordinate based networks are used to model the reflectance properties. Incident lighting is modelled as an implicit network with integrated directional embeddings [46]. We propose a differentiable rendering framework that takes as input the surface, reflectance parameters and illumination and renders polarization images under novel views. Given a set of multi-view polarization images, we jointly optimize parameters of the surface, reflectance parameters and incident illumination to minimize rendering loss.

Contributions. Our contributions are as follows:

- **Polarized neural rendering:** We propose a framework to render polarization images from implicit representations of the object geometry, surface reflectance and illumination.
- **3D surface reconstruction:** Equipped with implicit surface representation and polarization cues, PANDORA outputs high quality surface normal, signed distance field and mesh representations of the scene.
- **Diffuse-specular separation:** We demonstrate accurate diffuse-specular radiance decomposition on real world scenes under unknown illumination.
- **Incident illumination estimation:** We show that PANDORA can estimate the illumination incident on the object with high fidelity.

Assumptions. In this work, we assume that the incident illumination is completely unpolarized. The object is assumed to be opaque and to be made up of dielectric materials such as plastics, ceramics etc as our polarimetric reflectance model doesn’t handle metals. We focus on direct illumination light paths. Indirect illumination and self-occlusions are currently neglected.

2 Related Work

Inverse Rendering . The goal of inverse rendering is to recover scene parameters from a set of associated images. Inverse rendering approaches traditionally rely on multi-view geometry [40, 39], photometry [5] and structured lighting [31, 34] for 3D reconstruction [47], reflectance separation [31, 23], material characterization [16] and illumination estimation [38, 12]. Due to the ill-posed nature of inverse rendering, these approaches often require simplifying assumptions on the scene such as textured surfaces, Lambertian reflectance, direct illumination and simple geometry. Methods that aim to work in generalized scene settings involve incorporating scene priors [18, 54, 9], iterative scene optimization using differentiable rendering [22, 55] and exploiting different properties of light such as polarization [57], time-of-flight [53] and spectrum [21].

Neural Inverse Rendering . Recent emergence of neural implicit representations [49] has led to an explosion of interest in neural inverse rendering [44]. Neural implicit representations use a coordinate-based neural network to represent a visual signals such as images, videos, and 3D objects [41, 35, 27]. These representations are powerful because the resolution of the underlying signal is limited only by the network capacity, rather than the discretization of the signal. Interest from the vision community originated largely due to neural radiance field (NeRF) [28], which showed that modelling radiance using implicit representations leads to high-quality novel view synthesis.

Since the advent of NeRF, several works have exploited neural implicit representations for inverse rendering applications. IDR [51], UNISURF [33], NeuS [48] and VolSDF [52] demonstrate state-of-the-art 3D surface reconstruction from multi-view images by extending NeRF’s volume rendering framework to handle implicit surface representations. Accurate surface normals are crucial for modelling polarization and reflectance. Thus, we use ideas from one such work, VolSDF [52], as a build block in PANDORA.

NeRF models the net outgoing radiance from a scene point in which both the material properties and the lighting are mixed. Several approaches such as NeRV [43], NeRD [7], NeuralPIL [8], PhySG [56], RefNeRF [46] have looked at decomposing this radiance into reflectance components and illumination. PhySG and NeuralPIL employ spherical Gaussian and data-driven embeddings to model the scene’s illumination and reflectance. RefNeRF introduces integrated directional embeddings (IDEs) to model radiance from specular reflections and illumination and demonstrates improved novel view synthesis. Inspired from RefNeRF, we incorporate IDEs in our framework. Equipped with IDEs, implicit surface representation and polarimetric acquisition, PANDORA demonstrate better radiance decomposition than the state-of-the-art techniques, NeuralPIL and PhySG (Fig. 4,5, Table 1)

Polarimetric Inverse Rendering . Polarization strongly depends on the surface geometry leading to several single view depth and surface normal imaging approaches [29, 42, 3, 4, 20]. Inclusion of polarization cues has also led to enhancements in multi-view stereo [11, 57, 15, 14], SLAM [50] and time-of-flight imaging [17]. The diffuse and specular components of reflectance have distinct polarization properties and this distinction has been utilized for reflectance separation [26, 11, 19], reflection removal [25] and spatially varying BRDF estimation [13]. PANDORA exploits these polarimetric cues for 3D reconstruction, diffuse-specular separation and illumination estimation.

Traditionally acquiring polarization information required capturing multiple measurements by rotating a polarizer in front of the camera, unfortunately prohibiting fast acquisition. The advent of single-shot polarization sensors, such as the Sony IMX250MZR (monochrome) and IMX250MYR (color) [1] in commercial-grade off-the-shelf machine vision cameras has made polarimetric acquisition faster and more practical. These sensors have a grid of polarizers oriented at different angles attached on the CMOS sensor enabling the capture of polarimetric cues at the expense of spatial resolution. Various techniques have been proposed for

polarization and color demosaicking of the raw sensor measurements [30, 37]. In PANDORA we use such a camera FLIR BFS-U3-51S5P-C to capture polarization information for every view in a single shot.

3 Polarization as cue for radiance separation

Here we introduce our key insight on how polarimetric cues aide in decomposing radiance into the diffuse and specular components. First we derive the polarimetric properties of diffuse and specular reflectance. Then we demonstrate how these cues can aide in separating the combined reflectance by the means of a simple scene.

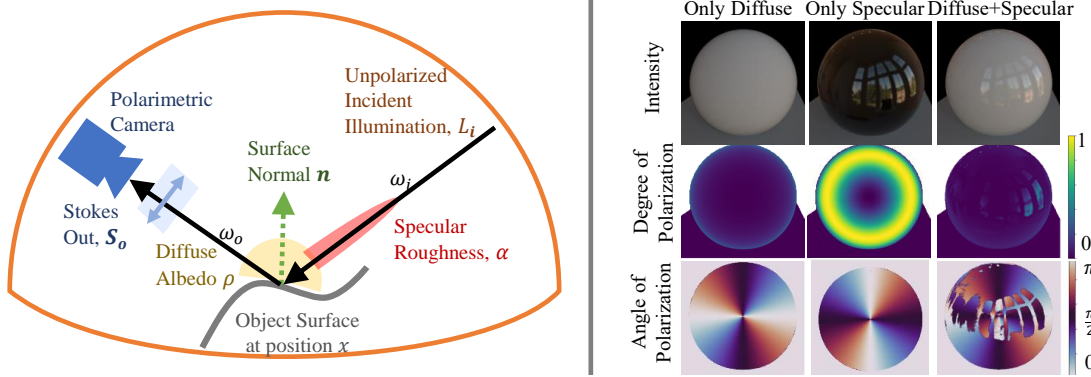


Figure 2: **Polarization as a cue for radiance decomposition** Left: Illustration of notations for our polarized image formation model. Right: Polarimetric cues for different radiance components. Diffuse radiance has a lower degree of polarization than the specular radiance. The polarization angles of diffuse and specular components are orthogonal. These polarization cues are used for better radiance decomposition.

3.1 Theory of polarized reflectance

Stokes Vector. The polarization state of light ray at \mathbf{x} along direction ω is modelled as Stokes vector comprising of four components, $S(\mathbf{x}, \omega) = [S_0 \ S_1 \ S_2 \ S_3]$ [10]. We assume there is no circular polarization and thus neglect S_3 . The Stokes vector can be parametrized as a function of three polarimetric cues: the total intensity, $L_o = S_0$, degree of polarization, $\beta_o = \sqrt{S_1^2 + S_2^2}/S_0$ and angle of polarization, $\phi_o = \tan^{-1}(S_2/S_1)/2$.

Mueller Matrix. Upon interaction with an object’s surface, the polarization state of light changes. The Stokes vector after the interaction can be modelled as the matrix multiplication of the input Stokes vector with a 4×4 matrix, known as the Mueller matrix.

Polarimetric BRDF (pBRDF) model. The interaction of object for diffuse and specular components of the reflectance are different. The diffuse reflectance involves sub-surface scattering into the surface and then transmission out of the surface. The specular component can be modelled as a direct reflection from specular microfacets on the surface. The pBRDF model [4] model these interactions as Mueller matrices for the diffuse and specular polarized reflectance, which we denoted as H_d and H_s respectively.

Incident Stokes vector Considering illumination is from far away sources, the dependance of S_i on \mathbf{x} can be dropped. Assuming the polarization to be completely unpolarized:

$$S_i(\mathbf{x}, \omega_i) = L_i(\mathbf{x}, \omega_i)[1 \ 0 \ 0]^T, \quad (1)$$

Exitant Stokes vector From the pBRDF model, the output Stokes vector at every point can be decomposed into the matrix multiplication of diffuse and specular Mueller matrices, H_d and H_s , with the illumination Stokes vector S_i ,

$$S_o(\mathbf{x}, \omega_i) = \int_{\Omega} H_d \cdot S_i(\mathbf{x}, \omega_i) d\omega + \int_{\Omega} H_s \cdot S_i(\mathbf{x}, \omega_i) d\omega \quad (2)$$

From S_i and the pBRDF model, we derive that the outgoing Stokes vector at every point depends on the diffuse radiance L_d , specular reflectance f_s and the incident illumination L_i as,

$$S_o(\mathbf{x}, \omega_i) = L_d \begin{bmatrix} 1 \\ \beta_d(\theta_n) \cos(2\phi_n) \\ -\beta_d(\theta_n) \sin(2\phi_n) \end{bmatrix} + L_s \begin{bmatrix} 1 \\ \beta_s(\theta_n) \cos(2\phi_n) \\ -\beta_s(\theta_n) \sin(2\phi_n) \end{bmatrix}, \quad (3)$$

where we terms β_d/β_s depend on Fresnel transmission/reflection coefficients for the polarization components parallel and perpendicular to the plane of incidence, $T^{\parallel}/R^{\parallel}$ and T^{\perp}/R^{\perp}

$$\beta_d = \frac{T^{\perp} - T^{\parallel}}{T^{\perp} + T^{\parallel}}, \quad \beta_s = \frac{R^{\perp} - R^{\parallel}}{R^{\perp} + R^{\parallel}}, \quad (4)$$

The Fresnel coefficients, T/R solely depend on the elevation angle of the viewing ray with respect to the surface normal, $\theta_n = \cos^{-1}(\mathbf{n} \cdot \omega_o)$. ϕ_o denotes the azimuth angle of the viewing ray with respect to the surface normals, $\phi_n = \cos^{-1}(\mathbf{n}_o, \mathbf{y}_o)$, where \mathbf{n}_o is the normal vector perpendicular to the viewing ray and \mathbf{y}_o is the y axis of camera coordinate system. Please refer to Appendix A for the detailed derivation and functional forms for Fresnel coefficients.

Next we show how these polarimetric cues depend on the diffuse and specular reflectance and aid in radiance decomposition.

3.2 Polarimetric properties of diffuse and specular radiance

From Eq. 3, the polarimetric cues of the captured Stokes vector are

$$L_o = L_d + L_s, \quad \beta_o = L_d \beta_d + L_s \beta_s, \quad \phi_o = \tan^{-1}(-\tan(2\phi_n))/2$$

Fig. 2 shows polarimetric cues for a sphere scene for different reflectance properties. For only diffuse case (left), the degree of polarization increases with elevation angle and the angle of polarization is equal to the azimuth angle. For only specular case (middle), the degree of polarization increases with elevation angle until the Brewster's angle after which it reduces. The angle of polarization is shifted from azimuth angle by 90° . When both diffuse and specular reflectance are present (right), the polarimetric cues indicate if a region is dominated by diffuse or specular radiance. The specular areas have higher degree of polarization than diffuse areas. The two components have orthogonal angle of polarization.

4 Our Approach

We aim to recover the object shape, diffuse and specular radiance along with the incident illumination from multi-view images captured from a consumer-grade snapshot polarization camera. Fig. 3 summarizes our pipeline.

4.1 Input

PANDORA relies on the following inputs to perform radiance decomposition: 1) *Polarization Images*. We capture multiple views around the object with a 4 MP snapshot polarization camera [1] (Fig. 3(a)). These cameras comprise of polarization and Bayer filter arrays on the sensor to simultaneously capture color images for four different polarizer orientations at the expense of spatial resolution. We employ the demosaicking and post-processing techniques utilized in [57] to convert the raw sensor measurements into 4 MP RGB Stokes

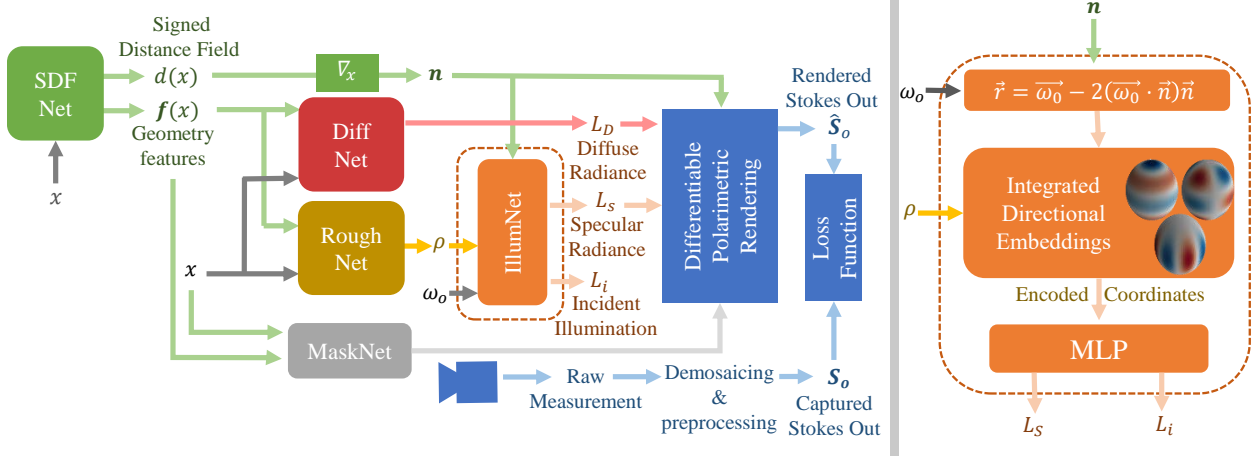


Figure 3: **PANDORA Pipeline**: Left: Our pipeline. We use coordinate-based networks to estimate surface normals, diffuse and specular radiance and incident illumination. From these parameters, we render exitant Stokes vector that is compared with captured Stokes vector and the loss is backpropagated to train the networks. Right: Detailed schematic of the Illumination Net

vector images. 2) *Camera poses*. We use COLMAP Structure-from-motion technique [40],[39] to calibrate the camera pose from the intensity measurements of the polarization images. Thus for any pixel in the captured images, the camera position \mathbf{o} and camera ray direction \mathbf{d} are known. An optional binary mask can also be used to remove signal contamination from the background. To create masks for real-world data, We use an existing object segmentation approach [36] for creating the object masks. The binary mask values are denoted as, $M(\mathbf{o}, \mathbf{d})$.

4.2 Implicit Surface estimation

The Stokes vector measured by camera ray given by \mathbf{o} and \mathbf{d} , the ray is sampled at \mathcal{T} points. For a sample on the ray with travel distance t , its location is denoted at $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$. The Stokes vector contribution of this sample depends on the scene opacity, $\sigma(\mathbf{r}(t))$ and exitant Stokes vector $S_o(\mathbf{r}(t), \mathbf{d})$. The observed Stokes vector, $S(\mathbf{o}, \mathbf{d})$ is denoted by the integral,

$$S(\mathbf{o}, \mathbf{d}) = \int_0^\infty T(t)\sigma(\mathbf{r}(t))S_o(\mathbf{r}(t), \mathbf{d})dt, \quad (5)$$

where $T(t) = \exp(-\int_0^\infty \sigma(\mathbf{r}(t)))$ is the probability that the ray travels to t without getting occluded.

For rendering surfaces, the ideal opacity should have a sharp discontinuity at the ray surface intersection. Thus accurately sampling \mathcal{T} points and consequently reconstructing sharp surfaces is challenging. High quality surface estimation is crucial for our approach as the polarization cues depend on the surface normals, (Eq. 3).

VolSDF [52] has demonstrated significant improvements in surface estimation by modelling the signed distance field d with a coordinate-based neural network. The opacity is then estimated as $\sigma(x) = \alpha\psi_\beta(-d(x))$ where α, β are learnable parameters and Ψ is the Cumulative Distribution Function of the Laplace distribution with zero mean and scale β . They also propose a better sampling algorithm for \mathcal{T} points utilizing the SDF representation. We follow the same algorithm as VolSDF for opacity generation. Similar to VolSDF, our pipeline comprises of an MLP, which we term SDFNet, that takes as input the position \mathbf{x} and outputs the signed distance field at that position \mathbf{x} along with geometry feature vectors \mathbf{f} useful for radiance estimation. The SDF model also provides us with surface normals, which are used in estimating specular radiance and polarimetric cues.

4.3 Neural Rendering Architecture

Diffuse Radiance Estimation. Diffuse radiance is invariant of the viewing direction and only depends on the spatial location. The geometry features from SDFNet and the position are passed through another coordinate-based MLP, denoted as DiffuseNet, to output the diffuse radiance at that position $L_D(\mathbf{x})$.

Specular Radiance Estimation. Unlike the diffuse radiance, the specular radiance depends on the viewing angle \mathbf{d} and the object roughness $\alpha\mathbf{x}$. First we estimate the object roughness using an coordinate-based MLP, RoughNet, similar in architecture to the DiffuseNet. For a certain object roughness, the obtained specular radiance involves integrating the specular BRDF along an incident direction factored by the incident illumination [6], which is a computationally expensive procedure that generally requires Monte Carlo. Inspired by [46], we instead use an IDE-based neural network to output the specular radiance, L_S from the estimated roughness, α and surface normals, \mathbf{n} . Moreover, on setting roughness close to zero, IllumNet also provides us the incident illumination, L_i .

Volumetric Masking We exploit object masks to ensure only the regions in the scene corresponding to the target object are used for radiance decomposition. Even when the background is zero, VolSDF estimates surface normals which have to be masked out to avoid incorrect quering of the IllumNet. Rather than using the 2D masks on the rendered images, we found that learning a 3D mask of the target object helps in training, especially in the initial iterations. This 3D mask $m(\mathbf{x})$ is 1 only for the positions \mathbf{x} that the object occupies and represent's the object's visual hull. We use this 3D mask to obtain the diffuse and specular radiance that is clipped to zero at background values,

$$L_D^m(\mathbf{x}) = m(\mathbf{x}) \cdot L_D(\mathbf{x}) \quad L_S^m(\mathbf{x}, \mathbf{d}) = m(\mathbf{x}) \cdot L_S(\mathbf{x}, \mathbf{d}) \quad (6)$$

The 3D mask is estimated using a coordinate-based MLP that we term MaskNet. This network is trained with the supervision of the input 2D object masks under different views. Similar to Eq. 5, the 3D mask values are accumulated along the ray and compared to the provided mask M using the binary cross entropy loss:

$$\mathcal{L}_{\text{mask}} = \mathbb{E}_{\mathbf{o}, \mathbf{d}} \text{BCE} \left(M(\mathbf{o}, \mathbf{d}), \hat{M}(\mathbf{o}, \mathbf{d}) \right), \quad (7)$$

where $\hat{M}(\mathbf{o}, \mathbf{d}) = \int_0^\infty T(t) \sigma(\mathbf{r}(t)) m(\mathbf{r}(t)) dt$.

Neural Polarimetric Rendering Using the masked diffuse L_D^m , masked specular L_S^m and the estimated surface normals \mathbf{n} , we can render the outgoing Stokes vector, $S_o(\mathbf{x}, \mathbf{d})$ from Eq. 3. On integrating outgoing Stokes vectors for points along the ray according to Eq. 5, we obtain the rendered Stokes vector $\hat{S}(\mathbf{x}, \mathbf{d})$.

4.4 Loss Function

We compare the rendered Stokes vector $\hat{S} = [\hat{s}_0, \hat{s}_1, \hat{s}_2]^T$ with the captured Stokes vector $S = [s_0, s_1, s_2]^T$ (§4.1) using an L1 loss. The loss is masked to remove the effect of background values. The s_1 and s_2 could have low values in regions having low degree of polarization (Fig. 2). We apply a weightage factor w_s on the loss for s_1 and s_2 outputs to further encourage the network to consider polarimetric cues in the training. The Stokes loss is modelled as:

$$\mathcal{L}_{\text{stokes}} = \mathbb{E}_{\mathbf{o}, \mathbf{d}} [M \cdot \|\hat{s}_0 - s_0\| + w_s \cdot M \cdot (\|\hat{s}_1 - s_1\| + \|\hat{s}_2 - s_2\|)] \quad (8)$$

Additionally, similar to VolSDF [52], we have the Eikonal loss, \mathcal{L}_{SDF} to encourage the SDFNet to approximate a signed distance field.

$$\mathcal{L}_{\text{SDF}} = \mathbb{E}_{\mathbf{o}, \mathbf{d}} (\|\mathbf{n}\| - 1)^2 \quad (9)$$

The net loss used to train all the networks described in the pipeline:

$$\mathcal{L}_{\text{net}} = \mathcal{L}_{\text{stokes}} + 0.1\mathcal{L}_{\text{SDF}} + \mathcal{L}_{\text{mask}} \quad (10)$$

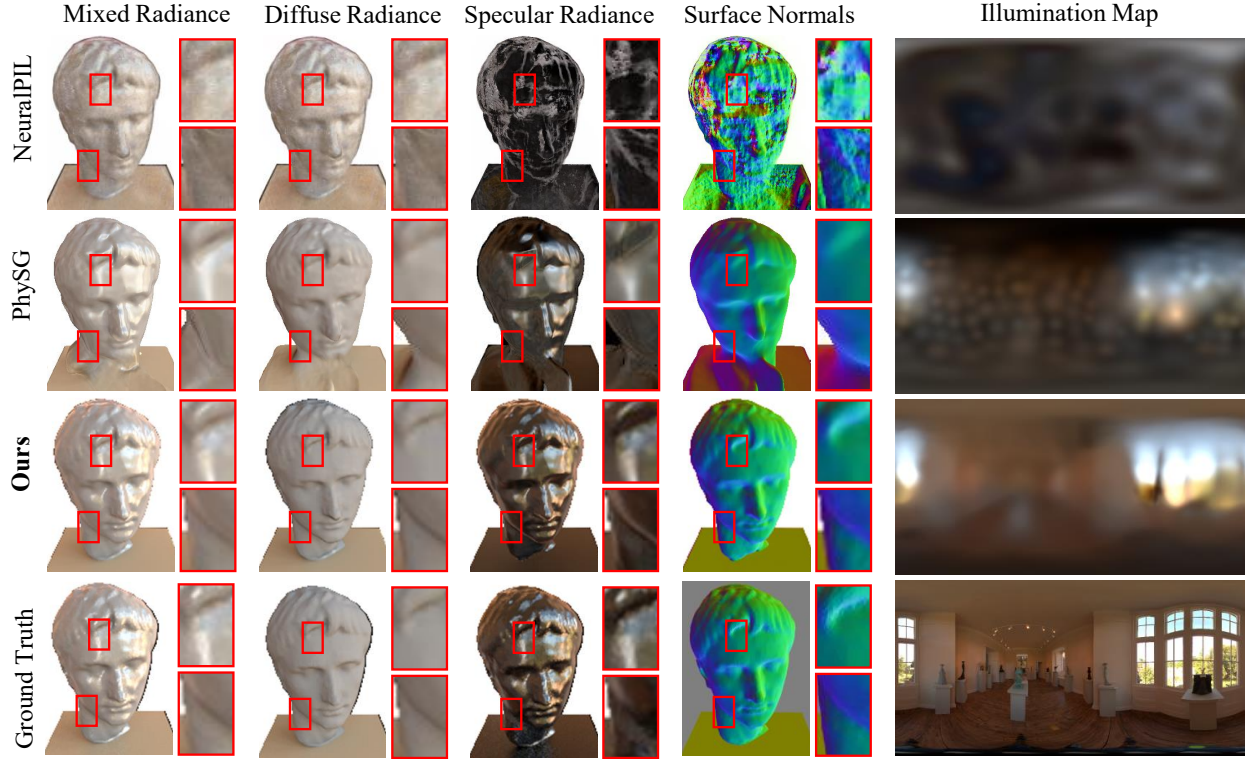


Figure 4: **Comparison of reflectance separation and surface normals with baselines on rendered dataset:** NeuralPIL fails to estimate correct normals and illumination on this challenging scene with strong specularities and 45 views. PhySG exhibits blurrier speculars and illumination along with artifacts in the reconstructed normals. PANDORA outputs sharp specularities, cleaner surface and more accurate illumination.

4.5 Implementation Details

All the networks are standard MLPs with 4 layers each. SDFNet has 256 hidden units per layer and the other networks have 512 hidden units. ReLU activations are used in intermediate layers. Final activation in DiffNet and MaskNet and the final activation in IllumNet and RoughNet is softplus. Please refer to Appendix B for additional implementation details of our framework.

5 Results and Evaluation

5.1 Datasets

We generate the following datasets for evaluating radiance decomposition.

1. Rendered Polarimetric Dataset (Fig. 4): Using Mitsuba2, we apply pBRDF on objects with complicated geometry and perform polarimetric rendering of multiple camera views under realistic environment lighting.
2. Real Polarimetric Dataset (Fig. 5,6,7): Using a snapshot polarimetric camera, we acquire multi-view polarized images of complex objects composed of materials with varying roughness, such as ceramics, glass, resin and plastic, under unstructured lighting conditions such as an office hallway. We also acquire the ground truth lighting using a chrome ball.

Please refer to Appendix B for additional details on the generation of these datasets and more examples from the datasets.



Scene	Approach	Diffuse		Specular		Mixed		Normals	Mesh
		PSNR ↑ (dB)	SSIM ↑	PSNR ↑ (dB)	SSIM ↑	PSNR ↑ (dB)	SSIM ↑	MAE ↓ (°)	HD ↓
Bust 	NeuralPIL	23.90	0.87	18.04	0.87	26.71	0.87	15.36	N/A
	PhySG	22.64	0.94	23.00	0.94	19.94	0.72	9.81	0.012
	Ours	25.82	0.81	22.96	0.75	22.79	0.79	3.91	0.003
Sphere 	NeuralPIL	13.09	0.55	12.92	0.55	20.04	0.66	38.73	N/A
	PhySG	21.76	0.76	18.90	0.76	17.93	0.70	8.42	0.011
	Ours	24.33	0.77	22.70	0.89	21.76	0.81	1.41	0.003

Table 1: **Quantitative evaluation on rendered scenes** We evaluate PANDORA and state-of-the-art methods on held-out testsets of 45 images for two rendered scenes. We report the peak average signal-to-noise ratio (PSNR) and structured similarity (SSIM) of diffuse, specular and net radiance, mean angular error (MAE) of surface normals and the Hausdorff distance (HD) of the reconstructed mesh. PANDORA consistently outperforms state-of-the-art in radiance separation and geometry estimation.

5.2 Comparisons with Baselines

We demonstrate that PANDORA excels in 3D reconstruction, diffuse-specular separation and illumination estimation compared to two existing state-of-the-art radiance decomposition baselines, NeuralPIL [8] and PhySG [56]. These baselines cannot exploit polarization cues and are run on radiance-only images using the public code implementations provided by the authors. We then show additional applications of PANDORA and an ablation study to analyse the crucial aspects of our algorithm.

3D Reconstruction The polarization cues directly depend on the surface normals (§3.2). Thus, inclusion of polarization cues, enhances multi-view 3D reconstruction. PANDORA reconstructs cleaner and more accurate surfaces such as jaw of the bust in Fig. 4 and the glass ball in Fig. 5. In table 1, we show that the mesh reconstructed by PANDORA has much lower Hausdorff distance with the ground truth mesh as compared to state-of-the-art. PANDORA also estimates more accurate surface normals as evaluated on a held-out test set.

Diffuse-Specular Separation The inherent ambiguity in separating diffuse and specular radiance components from intensity-only measurements leads to artefacts in existing techniques. For example, the black sphere in diffuse radiance reconstructed by NeuralPIL and PhySG contain faint specular highlights. Difference in polarization of diffuse and specular components enables PANDORA to obtain more accurate separation along with better combined radiance Fig. 4,5. In table 1, we show that PANDORA consistent outperforms state-of-the-art in peak signal-to-noise ratio (PSNR) and the structural similarity index measure (SSIM) of diffuse, specular and the net radiance images. We also provide the video of multi-view renderings from these diffuse, specular and mixed radiance fields that highlight the high quality of PANDORA’s separation.

Apart from the radiance, PANDORA can also separate the polarization properties of the object’s diffuse and specular components (Fig. 6). Here, we see predicted cues match with our physical intuition: the AoP is orthogonal for the diffuse and specular components, while the DoP is higher for the specular component.

Illumination estimation In addition to reflectance separation, our method can also estimate the illumination incident on the object. The rendered bust in Fig. 4 has blurry specular highlights that make illumination estimation challenging. We observe that NeuralPIL fails to estimate the correct lighting. PhySG employs spherical Gaussians that result in blurrier and more sparse reconstruction. PANDORA provides the best reconstruction with sharper walls and edges of the window.

Similarly, we also perform illumination estimation on real-world data (Fig. 7). We show results on data captured in two different environments. Fig. 7(left) is captured on a lab table with a long bright linear LED with dim ambient light. Fig. 7(right) is captured in a office hallway with many small tube-lights and bright walls. PhySG reconstruction is blurrier especially for the walls and comprises of color artifacts. PANDORA can recover high quality illumination that accurately matches the ground truth illumination as captured by replacing the object with a chrome ball.

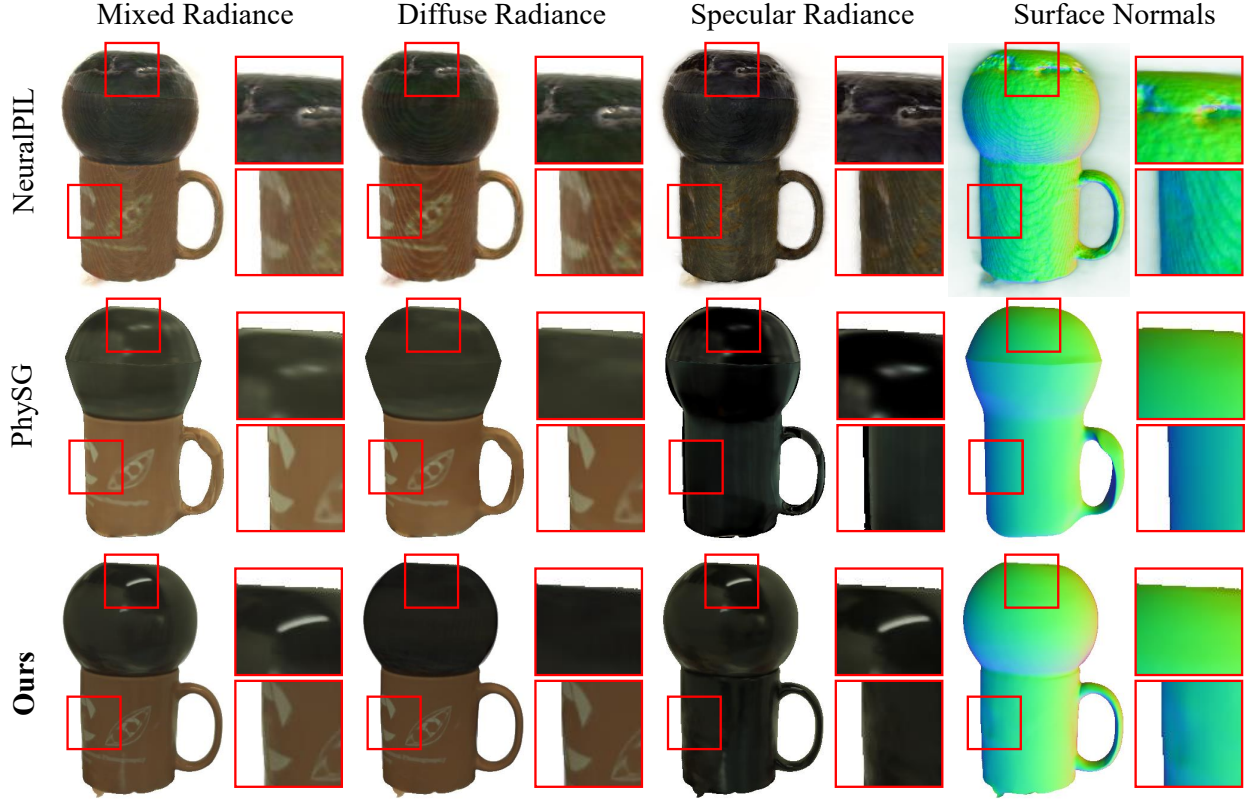


Figure 5: **Reflectance separation and surface normal estimation on real dataset** The decomposition using PhySG and NeuralPIL on intensity-only images has artifacts such as the specular highlights bleeding into the diffuse component and surface normals. PANDORA on polarized images produces accurate diffuse radiance, models the sharp specular reflections and reconstructs precise surface geometry.

5.3 Additional applications

The decomposed radiance field from PANDORA enables not only to render the object under novel views but also to change the object’s appearance under novel views by altering the separated diffuse and specular radiance fields. We demonstrate this application under Fig. 1(c). We perform polarimetric rendering from the learned PANDORA model under a new viewpoint. The rendered polarization (Fig. 1(c) bottom left) is consistent with the captured polarization. As PANDORA decomposes radiances, we can alter the diffuse component without affecting the specular reflections. For example, we assign pink albedo to the object by removing the G component of radiance without altering the color of the specularities in (Fig. 1(c) top left). To make the object look metallic, we render only the specular component with the Fresnel reflectance R^+ set to 1 (Fig 1(c) top right). To obtain rougher appearance (Fig 1(c) bottom right), we multiply the roughness parameter with a factor 3 before passing to the IllumNet. Please refer to project webpage for multi-view renderings of the separated reflectance fields and the changed appearance.

5.4 Ablation study: Role of polarization and IllumNet

Polarization and illumination modelling are key aspects of PANDORA. Here we analyse the role of these components by devising the following experiments

1. Ours w/o IllumNet w/o pol: We set Stokes loss weightage factor w_s (Eq. 8) as 0 to constrain PANDORA to just use S_o component, i.e., intensities for radiance decomposition. Also, instead of modelling illumination and roughness with neural networks, we directly model the specular radiance with a neural network same as the conventional RadianceNet in VolSDF.

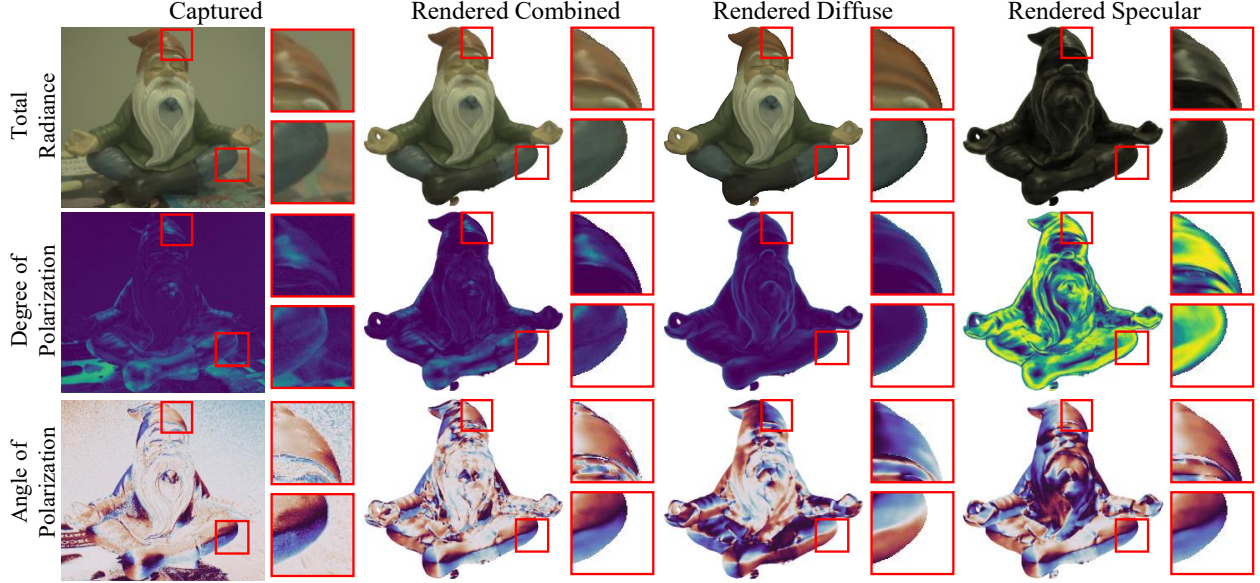


Figure 6: **Polarimetric diffuse-specular separation on real-world objects.** PANDORA can separate out diffuse and specular polarimetric properties from captured polarized images. As expected, rendered specular component has higher degree of polarization and polarization angle is orthogonal to the diffuse component. Polarization properties of the net rendered image match to that of the captured image.

2. Ours w/o pol: We set w_s as 0. But keep the IllumNet. So, this model has the same architecture as PANDORA. But it is trained on only the intensity.

We then train these two models and PANDORA on the same data with the same training scheme. As shown in Fig. 8, inclusion of the illumination modelling and polarization information significantly improves PANDORA’s performance. The model without IllumNet and polarization, exhibits strong artefacts of specular highlights in the diffuse and fails to capture the smaller specularities. Removing just polarization leads to worse illumination estimation, bleeding of diffuse into the specular and texture artefacts in the normals. Equipped with polarization information and correct illumination modelling, PANDORA outputs sharper diffuse texture, accurately handles the small specularities and even captures the subtle bumps on the object’s back.

6 Conclusion and Discussion

We have proposed PANDORA a novel neural inverse rendering algorithm that achieves state-of-the-art performance in reflectance separation and illumination estimation. PANDORA achieves this by using polarimetric cues and an SDF-based implicit surface representation. We have demonstrated the success of our approach on both simulated data that was generated by a physics based renderer, and real-world data captured with a polarization camera. Finally, we compared against similar approaches and demonstrated superior surface geometry reconstruction and illumination estimation. We believe PANDORA would pave the way for exciting ideas in the space of polarimetric and neural inverse rendering.

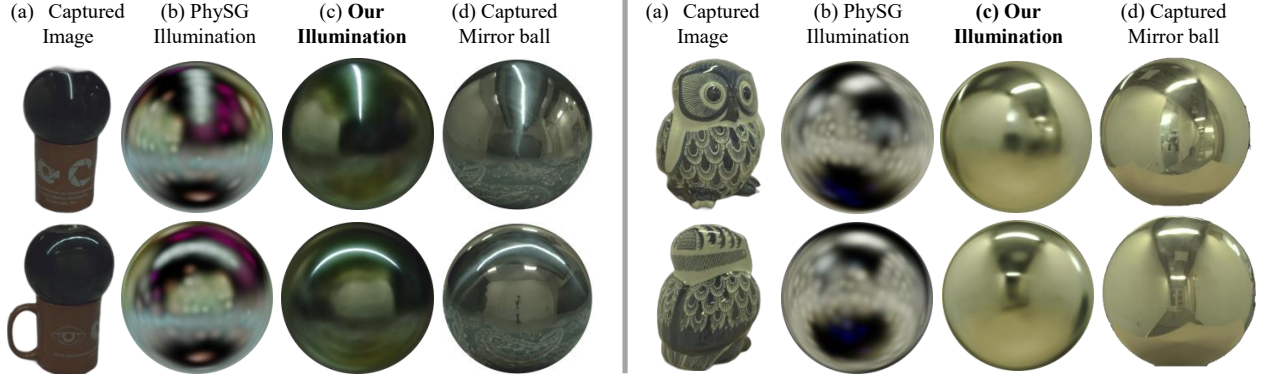


Figure 7: **Incident illumination estimated from real object** We visualize the illumination estimated on a mirror ball viewed from two train viewpoints. We also capture a mirror ball placed at a similar viewpoints. PhySG models the illumination using spherical Gaussians and leads to blurrier reconstruction with artefacts. PANDORA’s estimation has higher sharpness and accuracy.

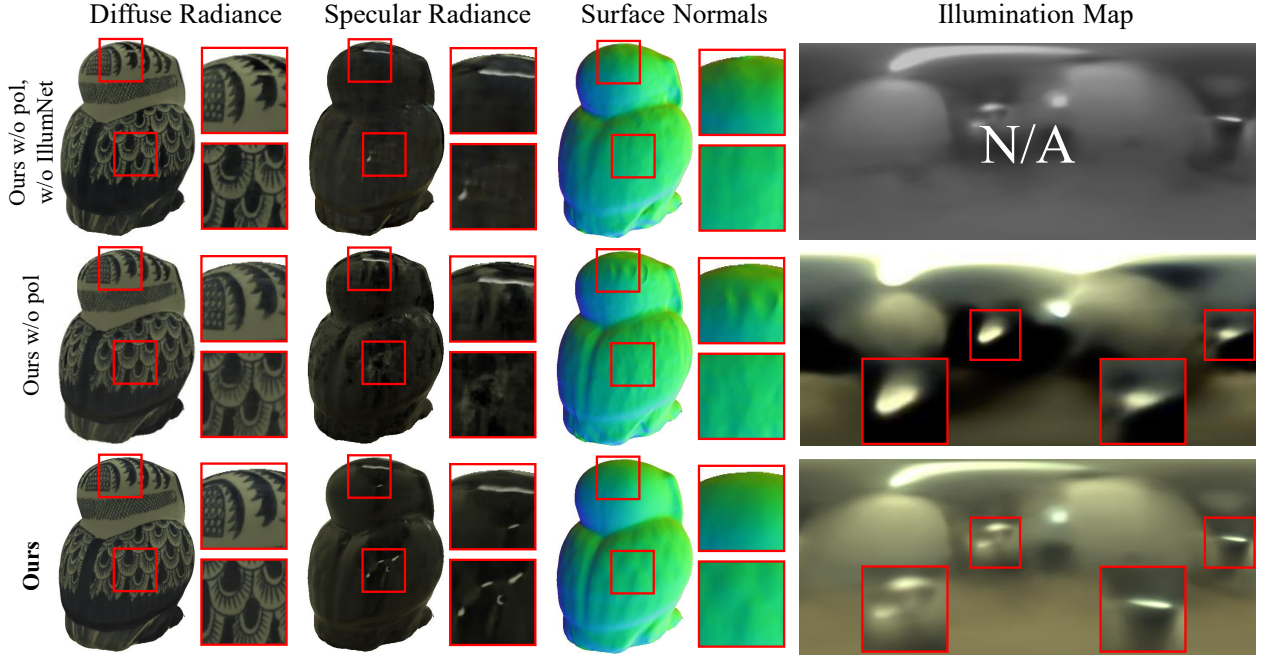


Figure 8: **Ablation study: Role of polarization and IllumNet** We devise two ablation experiments by running PANDORA on intensity-only images without IllumNet(top row) and with IllumNet(middle row). Without polarization and correct illumination modelling, there are texture artefacts in specular and surface normals due to the ambiguity in texture decomposition. Polarimetric cues and IllumNet help PANDORA in resolving such ambiguities resulting in finer quality reconstructions with sharper diffuse texture, cleaner surface normals and accurate lighting estimation.

References

- [1] Sony polarization image sensors. <https://www.sony-semicon.co.jp/e/products/IS/industry/product/polarization.html>, 2021. Accessed: 2021-09-25.
- [2] Alex Yu and Sara Fridovich-Keil, M. Tancik, Q. Chen, B. Recht, and A. Kanazawa. Plenoxels: Radiance fields without neural networks, 2021.
- [3] Y. Ba, A. Gilbert, F. Wang, J. Yang, R. Chen, Y. Wang, L. Yan, B. Shi, and A. Kadambi. Deep shape from polarization. 2020.
- [4] S.-H. Baek, D. S. Jeon, X. Tong, and M. H. Kim. Simultaneous acquisition of polarimetric svbrdf and normals. *ACM Trans. Graph.*, 37(6):268–1, 2018.
- [5] J. T. Barron and J. Malik. Shape, illumination, and reflectance from shading. *IEEE transactions on pattern analysis and machine intelligence*, 37(8):1670–1687, 2014.
- [6] J. T. Barron, B. Mildenhall, M. Tancik, P. Hedman, R. Martin-Brualla, and P. P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields, 2021.
- [7] M. Boss, R. Braun, V. Jampani, J. T. Barron, C. Liu, and H. P. Lensch. Nerf: Neural reflectance decomposition from image collections. In *IEEE International Conference on Computer Vision (ICCV)*, 2021.
- [8] M. Boss, V. Jampani, R. Braun, C. Liu, J. Barron, and H. Lensch. Neural-pil: Neural pre-integrated lighting for reflectance decomposition. *Advances in Neural Information Processing Systems*, 34, 2021.
- [9] Z. Chen, S. Nobuhara, and K. Nishino. Invertible neural brdf for object inverse rendering. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
- [10] E. Collett. *Polarized Light: Fundamentals and Applications*. CRC Press, United States, 1992.
- [11] Z. Cui, J. Gu, B. Shi, P. Tan, and J. Kautz. Polarimetric multi-view stereo. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 369–378, 2017. doi: 10.1109/CVPR.2017.47.
- [12] P. E. Debevec and J. Malik. Recovering high dynamic range radiance maps from photographs. In *Proceedings of the 24th Annual Conference on Computer Graphics and Interactive Techniques, SIGGRAPH '97*, page 369–378, USA, 1997. ACM Press/Addison-Wesley Publishing Co. ISBN 0897918967. doi: 10.1145/258734.258884. URL <https://doi.org/10.1145/258734.258884>.
- [13] V. Deschaintre, Y. Lin, and A. Ghosh. Deep polarization imaging for 3d shape and svbrdf acquisition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021.
- [14] Y. Ding, Y. Ji, M. Zhou, S. B. Kang, and J. Ye. Polarimetric helmholtz stereopsis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5037–5046, 2021.
- [15] Y. Fukao, R. Kawahara, S. Nobuhara, and K. Nishino. Polarimetric normal stereo. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 682–690, 2021. doi: 10.1109/CVPR46437.2021.00074.
- [16] I. Gkioulekas, S. Zhao, K. Bala, T. Zickler, and A. Levin. Inverse volume rendering with material dictionaries. *ACM Trans. Graph.*, 32(6), nov 2013. ISSN 0730-0301. doi: 10.1145/2508363.2508377. URL <https://doi.org/10.1145/2508363.2508377>.
- [17] A. Kadambi, V. Taamazyan, B. Shi, and R. Raskar. Polarized 3d: High-quality depth sensing with polarization cues. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3370–3378, 2015.
- [18] H. Kim, M. Zollöfer, A. Tewari, J. Thies, C. Richardt, and C. Theobalt. Inversefacenet: Deep single-shot inverse face rendering from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018.

- [19] J. Kim, S. Izadi, and A. Ghosh. Single-shot layered reflectance separation using a polarized light field camera. 2016.
- [20] C. Lei, C. Qi, J. Xie, N. Fan, V. Koltun, and Q. Chen. Shape from polarization for complex scenes in the wild. *arXiv preprint arXiv:2112.11377*, 2021.
- [21] C. Li, Y. Manno, and M. Okutomi. Spectral mvir: Joint reconstruction of 3d shape and spectral reflectance. In *2021 IEEE International Conference on Computational Photography (ICCP)*, pages 1–12. IEEE, 2021.
- [22] T.-M. Li, M. Aittala, F. Durand, and J. Lehtinen. Differentiable monte carlo ray tracing through edge sampling. *ACM Trans. Graph. (Proc. SIGGRAPH Asia)*, 37(6):222:1–222:11, 2018.
- [23] S. Lin, Y. Li, S. B. Kang, X. Tong, and H.-Y. Shum. Diffuse-specular separation and depth recovery from image sequences. In *European Conference on Computer Vision (ECCV)*, pages 210–224, May 2002.
- [24] Lucidrains. Se3 transformer - pytorch. <https://github.com/lucidrains/se3-transformer-pytorch>, 2021.
- [25] Y. Lyu, Z. Cui, S. Li, M. Pollefeys, and B. Shi. Reflection separation using a pair of unpolarized and polarized images. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/d47bf0af618a3523a226ed7cada85ce3-Paper.pdf>.
- [26] W.-C. Ma, T. Hawkins, P. Peers, C.-F. Chabert, M. Weiss, P. E. Debevec, et al. Rapid acquisition of specular and diffuse normal maps from polarized spherical gradient illumination. *Rendering Techniques*, 2007(9):10, 2007.
- [27] L. Mescheder, M. Oechsle, M. Niemeyer, S. Nowozin, and A. Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [28] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.
- [29] Miyazaki, Tan, Hara, and Ikeuchi. Polarization-based inverse rendering from a single view. In *Proceedings Ninth IEEE International Conference on Computer Vision*, pages 982–987 vol.2, 2003. doi: 10.1109/ICCV.2003.1238455.
- [30] M. Morimatsu, Y. Monno, M. Tanaka, and M. Okutomi. Monochrome and color polarization demosaicking using edge-aware residual interpolation. In *2020 IEEE International Conference on Image Processing (ICIP)*, pages 2571–2575. IEEE, 2020.
- [31] S. K. Nayar, G. Krishnan, M. D. Grossberg, and R. Raskar. Fast separation of direct and global components of a scene using high frequency illumination. *ACM Trans. Graph.*, 25(3):935–944, jul 2006. ISSN 0730-0301. doi: 10.1145/1141911.1141977. URL <https://doi.org/10.1145/1141911.1141977>.
- [32] M. Nimier-David, D. Vicini, T. Zeltner, and W. Jakob. Mitsuba 2: A retargetable forward and inverse renderer. *Transactions on Graphics (Proceedings of SIGGRAPH Asia)*, 38(6), Dec. 2019. doi: 10.1145/3355089.3356498.
- [33] M. Oechsle, S. Peng, and A. Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5589–5599, 2021.
- [34] M. OToole, J. Mather, and K. N. Kutulakos. 3d shape and indirect appearance by structured light transport. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 38(07):1298–1312, jul 2016. ISSN 1939-3539. doi: 10.1109/TPAMI.2016.2545662.

- [35] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [36] X. Qin, Z. Zhang, C. Huang, M. Dehghan, O. R. Zaiane, and M. Jagersand. U2-net: Going deeper with nested u-structure for salient object detection. *Pattern Recognition*, 106:107404, 2020.
- [37] S. Qiu, Q. Fu, C. Wang, and W. Heidrich. Linear polarization demosaicking for monochrome and colour polarization focal plane arrays. In *Computer Graphics Forum*. Wiley Online Library, 2021.
- [38] R. Ramamoorthi and P. Hanrahan. A signal-processing framework for inverse rendering. In *Proceedings of the 28th Annual Conference on Computer Graphics and Interactive Techniques, SIGGRAPH '01*, page 117–128, New York, NY, USA, 2001. Association for Computing Machinery. ISBN 158113374X. doi: 10.1145/383259.383271. URL <https://doi.org/10.1145/383259.383271>.
- [39] J. L. Schönberger and J.-M. Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [40] J. L. Schönberger, E. Zheng, M. Pollefeys, and J.-M. Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016.
- [41] V. Sitzmann, J. N. Martel, A. W. Bergman, D. B. Lindell, and G. Wetzstein. Implicit neural representations with periodic activation functions. In *arXiv*, 2020.
- [42] W. A. Smith, R. Ramamoorthi, and S. Tozza. Height-from-polarisation with unknown lighting or albedo. *IEEE transactions on pattern analysis and machine intelligence*, 41(12):2875–2888, 2018.
- [43] P. P. Srinivasan, B. Deng, X. Zhang, M. Tancik, B. Mildenhall, and J. T. Barron. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *CVPR*, 2021.
- [44] A. Tewari, J. Thies, B. Mildenhall, P. Srinivasan, E. Tretschk, Y. Wang, C. Lassner, V. Sitzmann, R. Martin-Brualla, S. Lombardi, et al. Advances in neural rendering. *arXiv preprint arXiv:2111.05849*, 2021.
- [45] turbosquid. Head_augustus 3d model, 2018. Accessed: 2022-01-23, generated by Cone of Vision. <https://www.turbosquid.com/3d-models/head-augustus-3d-model-1327693>.
- [46] D. Verbin, P. Hedman, B. Mildenhall, T. Zickler, J. T. Barron, and P. P. Srinivasan. Ref-nerf: Structured view-dependent appearance for neural radiance fields. *arXiv preprint arXiv:2112.03907*, 2021.
- [47] D. Vlasic, P. Peers, I. Baran, P. Debevec, J. Popović, S. Rusinkiewicz, and W. Matusik. Dynamic shape capture using multi-view photometric stereo. In *ACM SIGGRAPH Asia 2009 Papers, SIGGRAPH Asia '09*, New York, NY, USA, 2009. Association for Computing Machinery. ISBN 9781605588582. doi: 10.1145/1661412.1618520. URL <https://doi.org/10.1145/1661412.1618520>.
- [48] P. Wang, L. Liu, Y. Liu, C. Theobalt, T. Komura, and W. Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *NeurIPS*, 2021.
- [49] Y. Xie, T. Takikawa, S. Saito, O. Litany, S. Yan, N. Khan, F. Tombari, J. Tompkin, V. Sitzmann, and S. Sridhar. Neural fields in visual computing and beyond. *arXiv preprint arXiv:2111.11426*, 2021.
- [50] L. Yang, F. Tan, A. Li, Z. Cui, Y. Furukawa, and P. Tan. Polarimetric dense monocular slam. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3857–3866, 2018. doi: 10.1109/CVPR.2018.00406.
- [51] L. Yariv, Y. Kasten, D. Moran, M. Galun, M. Atzmon, B. Ronen, and Y. Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems*, 33:2492–2502, 2020.
- [52] L. Yariv, J. Gu, Y. Kasten, and Y. Lipman. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems*, 34, 2021.

- [53] S. Yi, D. Kim, K. Choi, A. Jarabo, D. Gutierrez, and M. H. Kim. Differentiable transient rendering. *ACM Transactions on Graphics (TOG)*, 40(6):1–11, 2021.
- [54] Y. Yu and W. A. Smith. Inverserendernet: Learning single image inverse rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [55] C. Zhang, L. Wu, C. Zheng, I. Gkioulekas, R. Ramamoorthi, and S. Zhao. A differential theory of radiative transfer. *ACM Trans. Graph.*, 38(6):227:1–227:16, 2019.
- [56] K. Zhang, F. Luan, Q. Wang, K. Bala, and N. Snavely. Physg: Inverse rendering with spherical gaussians for physics-based material editing and relighting. In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [57] J. Zhao, Y. Monno, and M. Okutomi. Polarimetric multi-view inverse rendering. 2020.

A Forward Model Derivation

In this section, we elaborate on the derivation of exitant Stokes vector as a function of diffuse and specular radiance as described in Eq. 3 of the main manuscript.

Diffuse Component In Eq. 2, we decompose the outgoing Stokes vector into diffuse and specular components. First we focus on the diffuse component. From the definition of H_d for pBRDF model [4] and the illumination Stokes vector defined in eq. 1, we obtain

$$H_d \cdot S_i = \rho(\mathbf{n} \cdot \mathbf{i}) L_i T_i^+ T_i^- \begin{bmatrix} T_o^+ \\ T_o^- \alpha_o \\ -T_o^- \delta_o \\ 0 \end{bmatrix}, \quad (11)$$

where ρ is the diffuse albedo, \mathbf{n} is the surface normal and \mathbf{i} is the incident illumination direction. With ϕ_n denoting the exitant azimuth angle w.r.t. the surface normal, we define α_o and δ_o as

$$\begin{aligned} \alpha_o &= \cos(2\phi_n) \\ \delta_o &= \sin(2\phi_n) \end{aligned} \quad (12)$$

We denote the term $\rho(\mathbf{n} \cdot \mathbf{i}) L_i T_i^+ T_i^-$ as the diffuse intensity L_D . The term $H_d \cdot S_i$ is independent of the viewing direction. Thus we obtain the first component of Eq.3

$$\int_{\Omega} H_d \cdot S_i(\mathbf{x}, \omega_i) d\omega = L_d \begin{bmatrix} 1 \\ T_o^- / T_o^+ \cos(2\phi_n) \\ -T_o^- / T_o^+ \sin(2\phi_n) \end{bmatrix} \quad (13)$$

Specular Component The specular exitant Stokes vector is obtained by substitution of H_s as defined in the pBRDF model [4] and S_i from eq. 1.

$$H_s \cdot S_i = L_i \frac{k_s DG}{4(\mathbf{n} \cdot \mathbf{o})} \begin{bmatrix} R^+ \\ R^- \chi_o \\ R^- \gamma_o \end{bmatrix}. \quad (14)$$

where k_s is the specular coefficient, \mathbf{o} is the exitant direction, D is the microfacet distribution and G is the microfacet shadowing term. With φ_h and φ_n denoting the incident and exitant azimuth angle w.r.t. the half angle \mathbf{h} respectively, we define χ_o and γ_o as

$$\begin{aligned} \chi_h &= \sin(2\varphi_h) \\ \gamma_h &= \cos(2\varphi_h) \end{aligned} \quad (15)$$

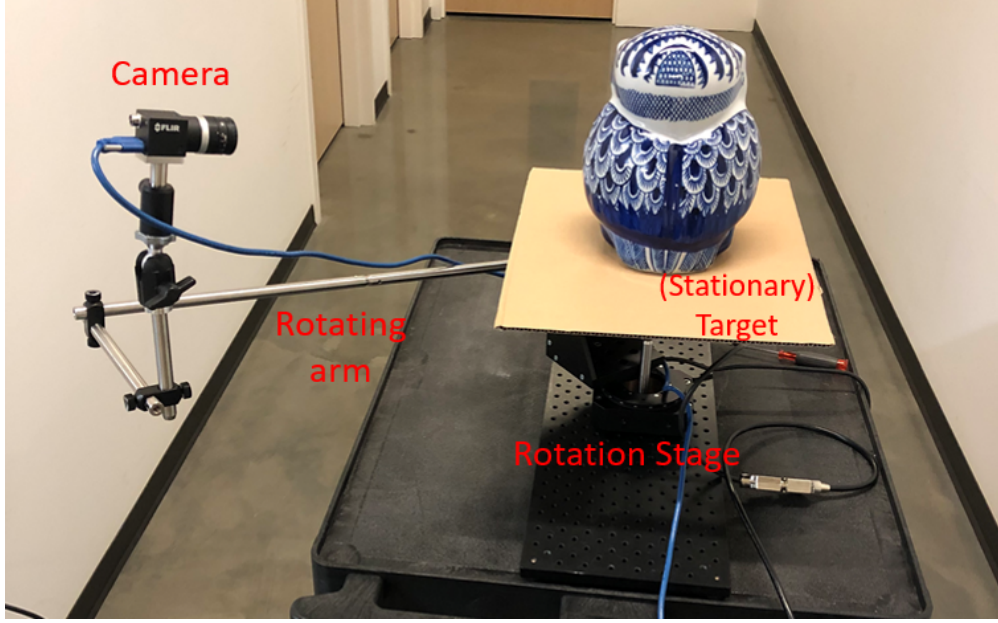


Figure 9: **Experimental Setup:** Above, is an image of our experimental setup. The target object is placed on the stationary section of a rotation stage, which is attached to an extended arm and the snapshot polarimetric camera. The camera capture polarimetric images from multiple angles under unstructured lighting while the target object remains still.

We denote $f_s = \frac{k_s DGR+}{4(\mathbf{n} \cdot \mathbf{o})}$. Theoretically χ and γ , depend on the half angles and not the geometric surface normals of the object. In practice, we observe that for realistic values of the roughness, χ and γ do not significantly deviate from the value obtained using surface normals instead of the half angle, i.e. $\chi_h \approx \sin(2\phi_h)$ and $\gamma_h \approx \cos(2\phi_n)$. As a result,

$$\int (H_s \cdot S_i) di = L_i \frac{k_s DG}{4(\mathbf{n} \cdot \mathbf{o})} \begin{bmatrix} R^+ \\ R^- \chi_o \\ R^- \gamma_o \end{bmatrix} \int f_s L_i di. \quad (16)$$

We denote $R^+ \int f_s L_i di$ as the specular radiance L_s and obtain the specular component of the output Stokes vector

B Implementation Details

Real world data was captured with a Blackfly S USB3 camera with Sony IMX250MYR Polarization-RGB sensor [1]. 35 images were captured for the Ball-Cup, Owl and Gnome objects under different lighting conditions as described in Table 2. The camera was placed along multiple angles distributed roughly equally along a circle around the target object using a portable setup as shown in Fig. 9. To capture the ground truth illumination map as shown in Fig. 12 last row, we use the same setup and flip the camera so that it points outside instead of the scene. Fish eye lens is used to increase the field of view and multi-view images are captured and stitched together to obtain the ground truth illumination map.

Rendering data generation Simulated data is generated using the Mitsuba2 renderer [32]. In Mitsuba2, we are able to set the material properties, camera angles, illumination, and imaging modality (polarized or unpolarized). We use a brdf that possesses equally weighted diffuse and dielectric (specular) components. We use 45 camera views distributed over all azimuth angles, and range from 25 to 50 degrees in elevation. Our

Camera Views

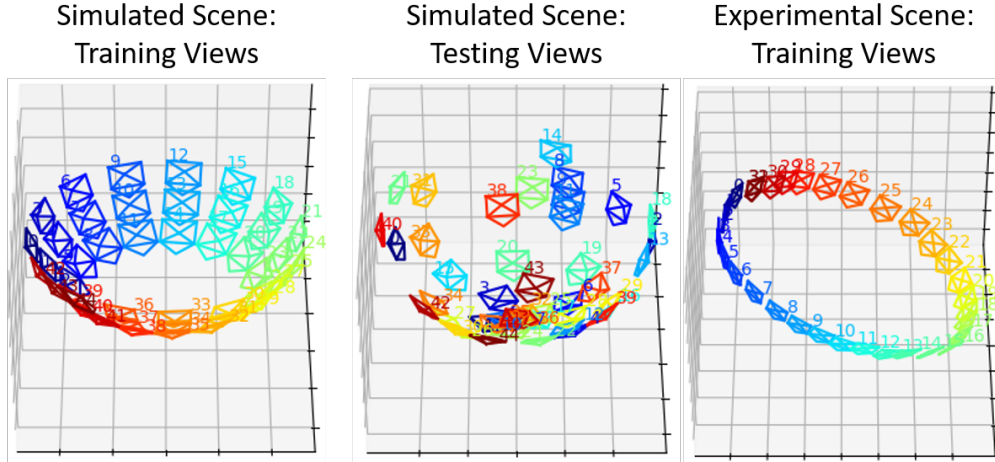


Figure 10: **Camera Views** Above, we show the camera positions for both the simulated and experimentally captured data.

two ground truth targets were a standard sphere and a bust shape obtained from [45]. The camera views are shown in Fig. 10.

Training details All training and testing was conducted on a server containing Nvidia 2080 Ti's. As stated in our main body, our DiffNet, MaskNet, RoughNet, and IllumNet were standard MLPs with 4 layers and a width of 512. Our SDFNet was an 8-layer MLP with a width of 256 and a single skip connection in the 4-th layer. Our training procedure uses several hyperparameters. The most relevant parameters include the weightage of the stokes vector loss, the weight of the mask network loss, the number of warm up iterations (before the stokes vector and specular components are estimated), and the total number of iterations. For real-world data we use 1000 warm-up iterations and 100,000 total iterations, while for simulated data we use 1500 warm-up iterations and 50,000 total iterations. We empirically found that a mask loss weightage and stokes loss weightage of 1.0 and 0.1, respectively, produced high-quality results. The diffuse and mask networks used a sigmoid activation function, while the specular and roughness networks used a softplus activation function to avoid vanishing gradients. Finally, for our SDFNet, MaskNet, RoughNet, and IllumNet, we used the frequency embeddings described by Mildenhall et al [28]. The frequencies of the embeddings were sampled in log-space from $2^0 - 2^6$ for the SDFNet and from $2^0 - 2^{10}$ for the MaskNet, RoughNet, and DiffNet. The integrated directional embeddings were used to embed the directional coordinates for the IllumNet, as described in more detail in the subsequent section.

Illumination Network Design The illumination network is responsible for calculating the incident illumination (the environment map) and the specular radiance, which is derived from Fresnel reflectance. To do this, the network accepts the reflected direction and the roughness as input. The roughness parameter is estimated by a separate network, while the reflected direction can be calculated from the predicted surface normals (using the geometry network) and the input viewing direction. Both inputs must be encoded through the IDE to help estimate the high frequency information and incorporate the effects of the roughness parameter, i.e. increase the blurring of the predicted illumination as the roughness gets larger. For the input to our IllumNet, we used degree $L \in \{1, 2, 4\}$ spherical harmonics with order $m \in [-L, L]$ for the IDE's.

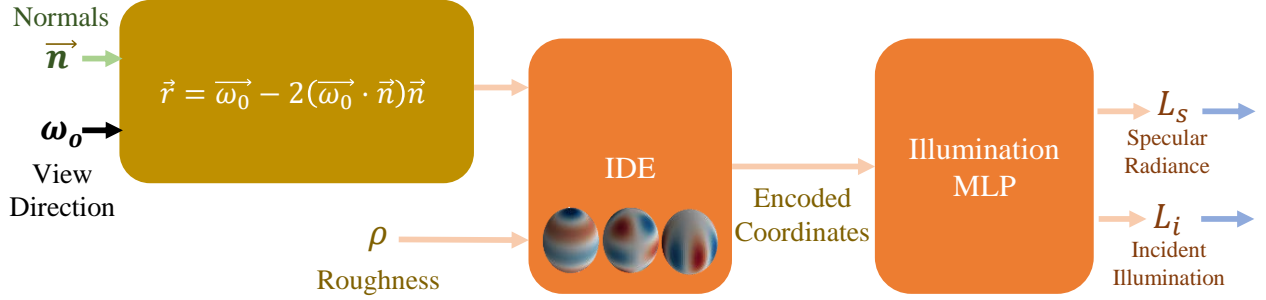


Figure 11: **Illumination Network Design:** The illumination network accepts the reflected direction vector and the predicted surface roughness as input. The reflected direction is calculated from the surface normal and viewing direction as shown above. The roughness and direction vector are encoded by the IDE before it is passed to the MLP which generates the predicted illumination and radiance based on fresnel reflectance.

C Additional Results



Scene	Approach	Diffuse		Specular		Mixed		Normals
		PSNR ↑ (dB)	SSIM ↑	PSNR ↑ (dB)	SSIM ↑	PSNR ↑ (dB)	SSIM ↑	MAE ↓ (°)
Bust 	NeuralPIL	23.90	0.87	18.04	0.87	26.71	0.87	15.36
	PhySG	22.64	0.94	23.00	0.94	19.94	0.72	9.81
	Ours no pol no Illum	28.29	0.968	21.13	0.906	22.29	0.951	7.89
	Ours no pol	25.78	0.956	18.23	0.856	22.50	0.927	4.83
	Ours	29.53	0.973	23.63	0.912	25.97	0.951	1.95
Sphere 	NeuralPIL	13.09	0.55	12.92	0.55	20.04	0.66	38.73
	PhySG	21.76	0.76	18.90	0.76	17.93	0.70	8.42
	Ours no pol no Illum	20.65	0.76	16.23	0.76	17.11	0.72	1.91
	Ours no pol	22.20	0.83	21.30	0.87	20.87	0.82	1.92
	Ours	24.29	0.84	21.29	0.88	21.29	0.83	1.04

Table 2: **Quantitative evaluation on rendered scenes** We evaluate PANDORA with state-of-the-art and ablation methods on held-out testsets of 45 images for two rendered scenes. We report the peak average signal-to-noise ratio (PSNR) and structured similarity (SSIM) of diffuse, specular and net radiance and mean angular error (MAE) of surface normals. PANDORA consistently outperforms state-of-the-art in radiance separation and geometry estimation.

In Fig. 12, we show additional qualitative comparisons with state-of-the-art inverse rendering technique, PhySG [56], and ablation model run on intensity-only images. In Fig. 13, we highlight the advantages of PANDORA over existing mesh optimization-based polarimetric inverse rendering technique, PMVIR [57]. We also report additional quantitative metrics on simulated and real data in Table 2 and Table 2 respectively. Please refer to the project webpage for videos showcasing our multi-view renderings.

D Analysis

Performance on out-of-distribution views As expected, for regions outside of the views in our training images, the estimation performs poorly. We see in Fig. 14 the network extrapolates a blob above the statue, in regions that are not heavily sampled during training. This affects our rendering when we sample rays in these regions (Fig. 14 panel 4). Finally, we see that by sampling rays only within a narrower region of interest, corresponding to locations with more training views, we obtain a correct estimate. We should note that in our main paper, the reported metrics do not account for this poor extrapolation as the images were rendered over a wider region of interest. So, the metrics were affected by artefacts in some of the rendered images shown in Fig. 14 panel 4. Metrics reported in Table 2 are with images rendered over smaller region of interest and do not have this artefact.

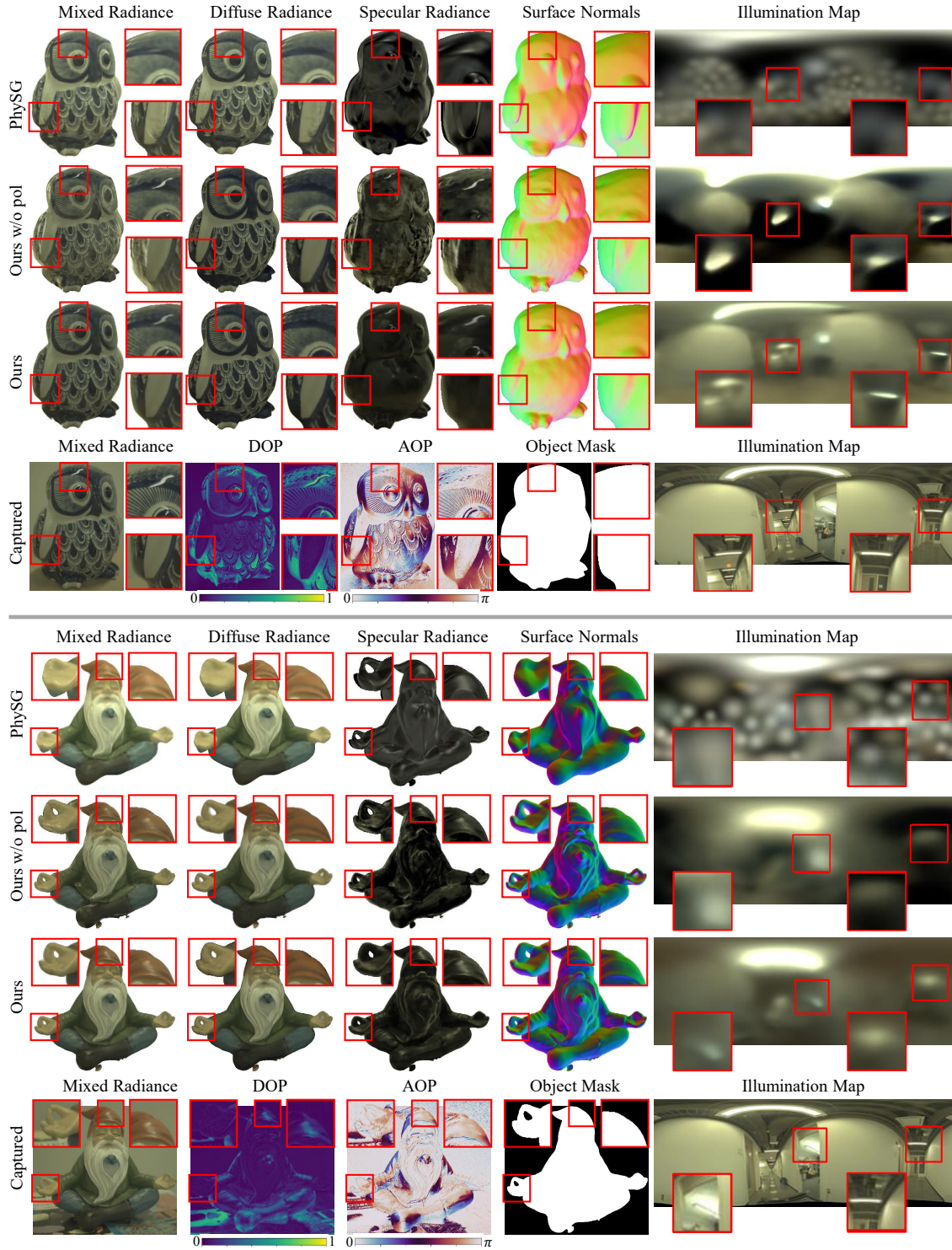


Figure 12: **Reflectance separation, surface normal reconstruction and illumination estimation on real dataset PANDORA** captures high frequency details in the surface normals and accurately models the specular highlights. Please view the project webpage for multi-view renderings of the same.





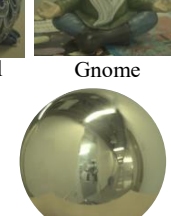
Scene												
	Ball-Cup	Owl	Gnome									
Lighting			Lab	Hallway	PhySG		Ours w/o pol		Ours			
					PSNR	SSIM	PSNR	SSIM	PSNR	SSIM		
					↑ (dB)	↑	↑ (dB)	↑	↑ (dB)	↑		
					Owl	Hallway	27.68	0.953	27.67	0.940	30.37	0.960
					Gnome	Hallway	30.31	0.986	28.42	0.984	29.15	0.984
					Ball-cup	Hallway	19.46	0.920	27.99	0.980	28.12	0.981
Ball-cup	Lab	14.00	0.950	23.52	0.953	26.92	0.970					

Table 3: **Quantitative evaluation on real scenes** We report the average PSNR and SSIM of the rendered intensity image over the training set for objects with different material properties and under different lighting conditions. PANDORA consistently outperforms PhySG and the ablation model that is devoid of the polarimetric cues.

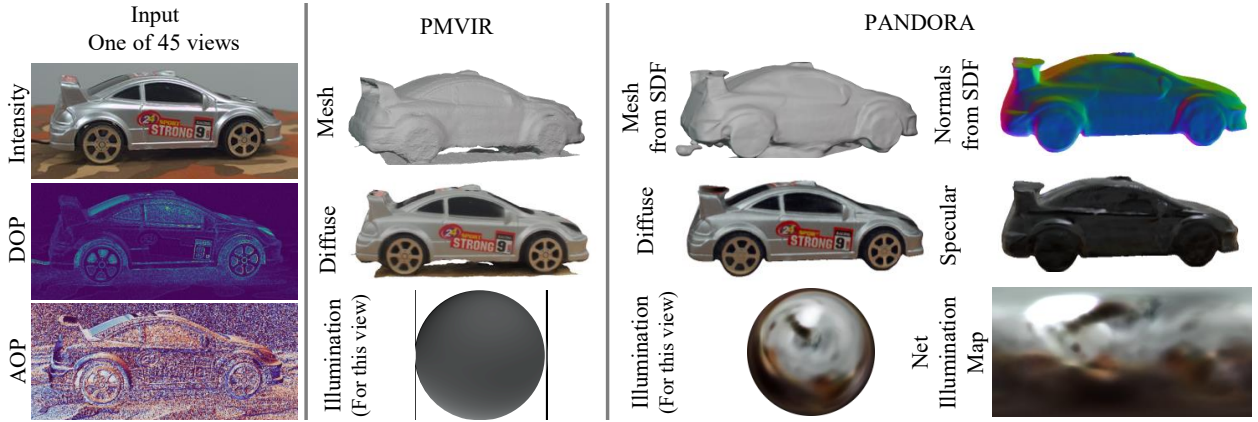


Figure 13: **Comparison with prior mesh-based polarimetric inverse rendering on real data** Utilizing similar multi-view snapshot polarimetric data as ours, PMVIR [57] recovers 3D mesh, diffuse color for mesh vertex and lighting based on diffuse shading. Neural implicit representations enable PANDORA to extract more from the same captured data. PANDORA learns the continuous signed distance field from which mesh and surface normals can be extracted. Apart from the diffuse color, PANDORA also outputs the specular radiance. Illumination estimated from PANDORA features sharp light source and the orange floor that better explain the captured data.

Effect of roughness on illumination estimation Above, we show the effects of the surface roughness on the estimated illumination map. As the surface roughness (α) increases, the associated, estimated environment map is increasingly blurred. The inset images show the ground truth specular reflection for each of the estimated environment maps. On the right-hand side, we show the associated spherical harmonic bases, which are used for the integrated directional encoding (IDE)¹ [46]. Recall that the IDE is used to encode the directional coordinates, which are passed as input to the illumination MLP. Increasing roughness decreases the impact of the higher frequency spherical harmonic bases, as shown on the right. This helps to supervise the desired blurring effect because the high-frequency components reduced.

Effect of roughness on polarimetric cues. In Fig. 16, we show using renderings from Mitsuba that the variation of polarimetric cues on varying roughness is less and the polarization of specular component is always distinct from the diffuse polarization under different levels of roughness.

¹The IDE visualization was generated using the ReF-NeRF implementation, with help in implementing the spherical harmonics transform from [2, 24]

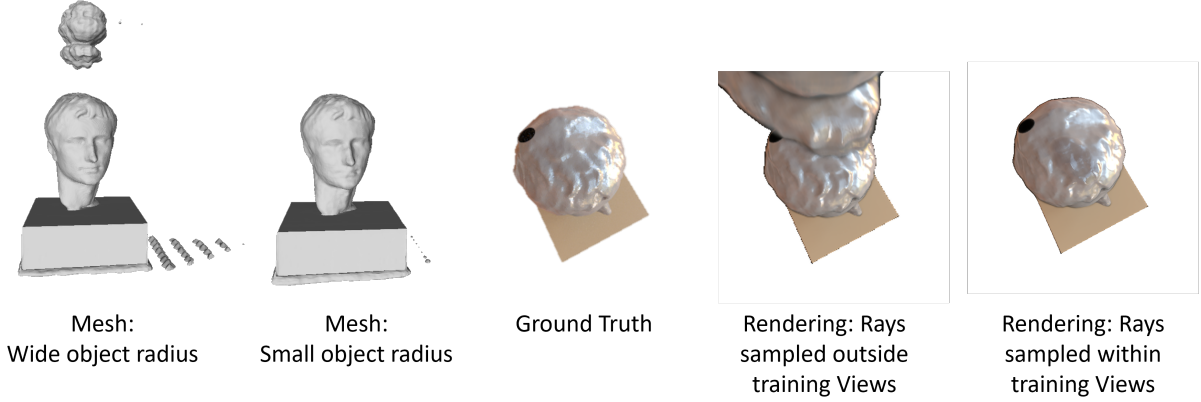


Figure 14: **Extrapolated Views Result:** We show the estimated mesh corresponding to regions that had lower sampled (panel 1). and higher sampled (panel 2) views. In addition, we show the resulting renderings when using more extrapolated rays (panel 4) versus without the extrapolated rays (panel 5).

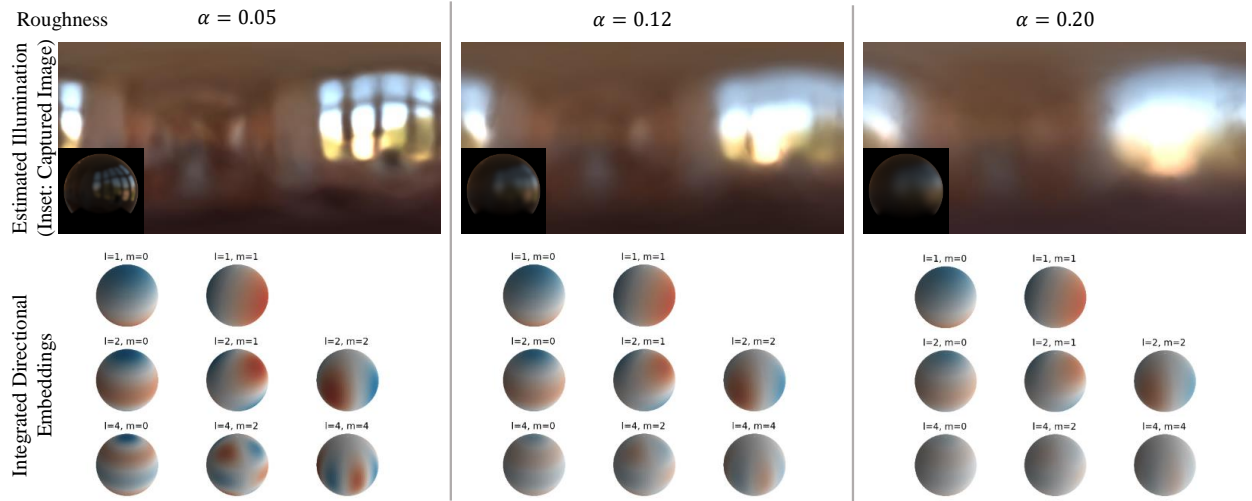


Figure 15: **Effect of roughness on illumination estimation:** Our illumination estimation accounts for the effects of surface roughness. As the roughness (parameterized by α) increases, there is an increasing blurring effect on the estimated environment map. The inset images shows the corresponding ground truth specular reflection as the surface roughness increases. The right side shows the effect of the increasing roughness on the spherical harmonic IDE's.

E Limitations

There are two main limitations to our current approach. Firstly, our method does not handle self-occlusions. This is more prominent in our simulated bust target, since the target geometry is not fully convex. We see dark patches in the estimated illumination map where the network cannot correctly estimate the illumination due to self-occlusions. In future work, this limitation may be resolved using a similar method as Verbin et al [46], in which a learnable “bottleneck” vector is used to model the target features that are not explained by other parts of the network.

Secondly, our method is not able to perform re-lighting. While PANDORA can perform diffuse-specular radiance separation, the incident illumination is baked into these radiances and it is challenging to estimate physically-based material properties, more specifically the material roughness and the diffuse albedo. While our network outputs an α parameter that tunes the rough appearance and models the effect of increasing roughness, such as blurred illumination map (Fig. 15), it does not truly estimate the physics-based roughness.

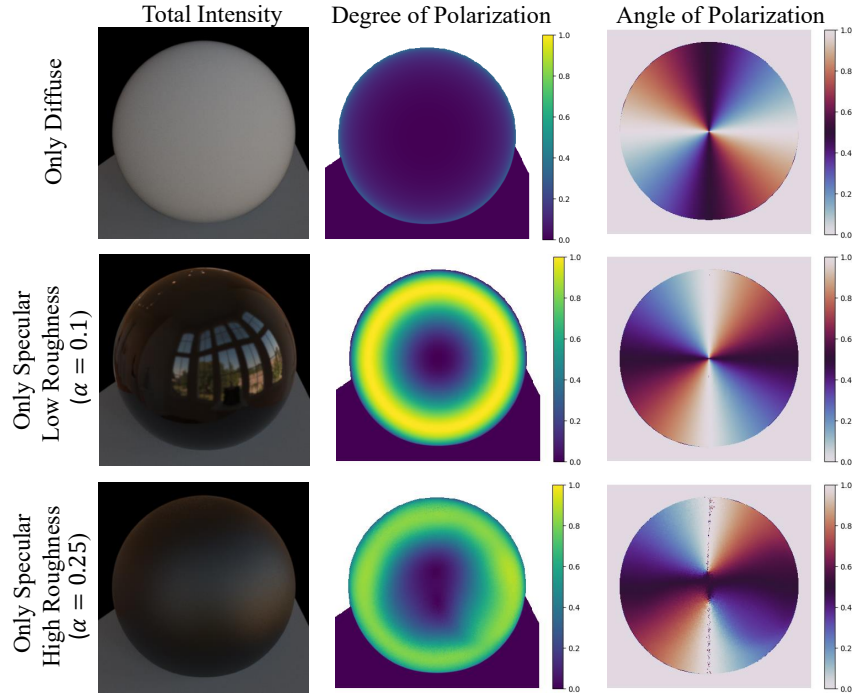


Figure 16: **Effect of specular roughness on polarimetric cues** We render polarimetric cues for a sphere object using the pBRDF model in Mitsuba2 with varying material properties. The variation of polarimetric cues is less under the realistic range of roughness. Our insight that the specular polarization is orthogonal in angle and higher in degree than the diffuse polarization remains applicable on varying specular roughness to realistic values. α denotes the roughness parameter of the Beckmann microfacet distribution.