# Surf-NeRF: Surface Regularised Neural Radiance Fields

Jack Naylor        Viorela Ila        Donald G. Dansereau

Australian Centre for Robotics
The University of Sydney
Darlington, NSW, 2006, Australia

{firstname.lastname}@sydney.edu.au

## Abstract

*Neural Radiance Fields (NeRFs) provide a high fidelity, continuous scene representation that can realistically represent complex behaviour of light. Despite recent works like Ref-NeRF improving geometry through physics-inspired models, the ability for a NeRF to overcome shape-radiance ambiguity and converge to a representation consistent with real geometry remains limited. We demonstrate how curriculum learning of a surface light field model helps a NeRF converge towards a more geometrically accurate scene representation. We introduce four additional regularisation terms to impose geometric smoothness, consistency of normals and a separation of Lambertian and specular appearance at geometry in the scene, conforming to physical models. Our approach yields improvements of 14.4% to normals on positionally encoded NeRFs and 9.2% on grid-based models compared to current reflection-based NeRF variants. This includes a separated view-dependent appearance, conditioning a NeRF to have a geometric representation consistent with the captured scene. We demonstrate compatibility of our method with existing NeRF variants, as a key step in enabling radiance-based representations for geometry critical applications. Project page: https://roboticimaging.org/Projects/SurfNeRF*

## 1. Introduction

Neural Radiance Fields (NeRFs) [27] provide an efficient coordinate based scene representation with wide applications to computer vision, robotics and beyond. The ray-based volumetric rendering formulation produces realistic novel views of a scene, encompassing view-dependent scene appearance. Complex appearances like specularity, transparency and interreflection lead to shape-radiance ambiguity where scene geometry is not uniquely represented, often producing imagined geometries.

Whilst re-parameterisation of the radiance in the scene as in Ref-NeRF [36] has shown improvements over previous
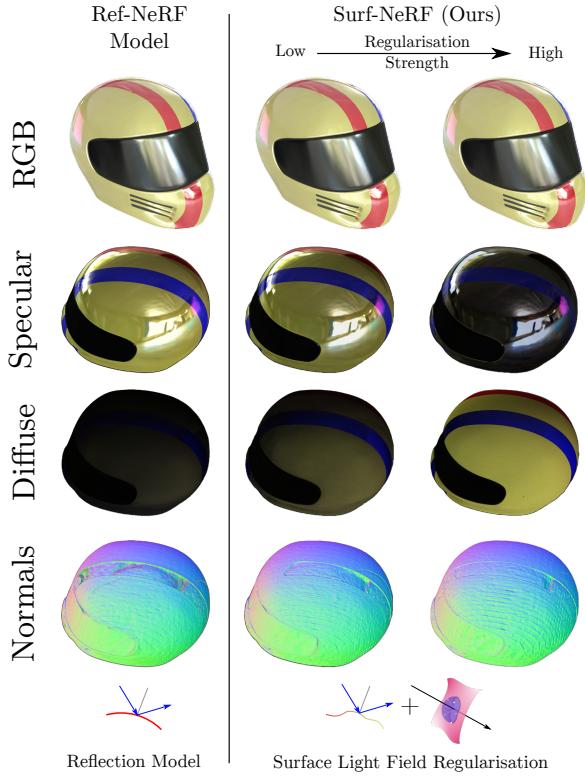


Figure 1. Surf-NeRF uses the properties of surface light fields to regularise a NeRF by querying the local region of samples. These samples are used to improve consistency and smoothness in these regions accumulating density at a surface, leading to improved geometry and more physically viable appearance components (ours, right) compared to current state-of-the-art ( [36], left) whilst maintaining visual fidelity. By changing how frequently regularisation occurs during training, we can control the strength of the effect.

state-of-the-art variants, there exists an incomplete separation of the scene's Lambertian appearance from the specular and no explicit constraint on the placement of density. It is still common for the NeRF to place regions of density behind surfaces or in front of the camera as detached geometry (*floaters*) to explain complex phenomena. This leads to a

1

non-realistic geometric scene representation, as in Figure 1.

Applications like robotics and 3D modelling require both accurate geometry and realistic appearance. The poor physical geometry currently poses a large hurdle to the widespread adoption of NeRFs as a scene representation for geometry critical tasks. Robotic manipulation and autonomous navigation and mapping require an accurate scene structure to minimise the gap between a representation and the real world, for example when grasping metal objects or navigating near reflective windows. Similarly, traditional structure-from-motion pipelines have great difficulty in reconstructing visually complex objects including reflective and transparent surfaces.

We address geometric inaccuracy arising from view-dependent phenomena with the insight that there are multiple formulations of the plenoptic function which can produce viable scene representations. A surface light field [41] describes the plenoptic function as light originating from a geometrically smooth surface. Our insight lies in that we can push the NeRF towards a more geometrically accurate representation in the same rendering framework using additional geometric and appearance based regularisation via curriculum learning. We use a first surface assumption to describe the location of geometry in the scene. Using a second sampling of the NeRF at these points, we regularise density to produce smoothly-varying normals and thin, continuous sheets of density which more realistically represent scene geometry. By enforcing view-dependent properties of a reflection light model at these points, we enable a separation of the Lambertian appearance of a scene from the view-dependent component, reducing shape-radiance ambiguity and encouraging more correct structure of the surface.

In this work, we make the following contributions:

1. We devise a novel regularisation approach which uses the structure of a neural radiance field to sample density, normals and appearance in the vicinity of geometry in the scene, allowing for additional representation-driven regularisation terms to be applied.

2. We apply local regularisation consistent with a surface light field radiance model, including geometric smoothness of density, local consistency of normals and a physically correct separation of Lambertian and specular appearance using a light interaction model.

3. We leverage curriculum learning of a NeRF towards a more accurate geometric scene representation which maintains visual fidelity whilst refining the density representation of the scene.

Whilst we benchmark our approach on state-of-the-art physics based NeRF variants, our methodology may also be applied to other NeRF frameworks.

This work is a key step in the deployment of NeRFs as a scene representation where both geometric and visual fidelity are critical, like robotic manipulation and navigation

in complex unstructured environments.

## 2. Related Work

**Neural Scene Representations:** Neural field approaches [14, 35] to scene representation produce continuous and often high visual fidelity depictions which are able to be queried anywhere within the training set, balancing visual and geometric fidelity. In this work, we seek to improve the geometry of a NeRF whilst maintaining visual fidelity.

NeRFs [27] leverage the efficiency of ray-based construction similar to light fields with a volumetric scene representation of points which emit light. Subsequent works have improved the fidelity of representation [3, 4, 6], recovery of a static scene [26, 32, 34] and its performance around view-dependent phenomena such as reflections [22, 36]. However, relying on a purely volumetric rendering approach still allows the network to imagine geometries to explain view-dependent phenomena, particularly in cases where regions of the scene are underconstrained [23, 29, 31] or where light does not follow the physical model of a straight ray through the scene [8]. Additional regularisation terms have been shown to dramatically improve the quality of rendering in these cases [23, 29]. Improving the accuracy of normals in the scene also significantly helps learning around complex appearance [25, 36]. More recent works [37] have shown that ray-tracing multiple reflection rays can improve the representation of reflections, albeit having extremely high computation cost. In this work we use properties of surface light fields to reduce the reliance of NeRFs on additional geometries. This approach explains complex visual phenomena entirely through view-dependent appearance, whilst maintaining the formulation of a single-pass volumetric rendering approach and its quality. We achieve this using a surface light field model and local characteristics of geometry and appearance of the scene.

**Surface Representations:** Signed distance fields (SDFs) have garnered significant attention [2, 16, 18, 30, 44, 45] as they can represent smooth and continuous geometries in space and may be generated from multi-view constraints alone. Applying view-dependent colour channels to the SDF [18, 38–40, 44] in a similar fashion to NeRFs has demonstrated improved accuracy and continuity of representation. For diffuse objects, this representation provides a smooth and high quality reconstruction, however, more highly view-dependent appearance and complex geometry (thin structures, concave geometries) are not well represented. Some works have included a reflection parameterisation [21], however struggle to reconstruct 3D scenes with fine detail, thin structures and sharp changes compared to volumetric scenes. Appearance and geometry have an intrinsic relationship within neural representations, with degradation in geometric accuracy often resulting in altered geometry to meet appearance [18]. Our work main-

tains a volumetric scene representation allowing for complex geometry, but with separated view-dependent and -independent appearance. Given NeRFs are a high fidelity visual representation, we are concerned only with introducing a surface-like structure to the density field through volumetric rendering. This improves the geometry and consistency of novel view rendering by considering the NeRF as learning density constrained to a surface with smooth view-dependent terms.

**Appearance Vs. Geometry:** Inverse rendering seeks to produce a definite [20] separation of the scenes appearance, geometry and environment, which are necessary to create realistic renderings of scene's under new conditions. Performing this without prior knowledge of the scene proves to be an immensely difficult task [11, 15, 43, 46], however, utilising physically-based rendering fused with learnt appearances [10,12] provides a sufficiently constrained framework to acquire the components of the scene. Radiance field approaches seeking to learn an accurate scene representation under few view scenarios have used depth consistency [33], additional depth supervision [49] and regularisation using priors [29]. Incorporating an understanding of the physical interactions of light within a scene and its effect on appearance [47] under a solid surface assumption enables accurate geometry to be learnt alongside appearance. In the presence of specularities and other visual phenomena it is difficult to disentangle where appearance has been baked into geometry [42]. Without the need for recovering a bidirectional reflectance distribution function (BRDF), light field approaches [19] provide a more generalised framework. This enables a clear separation between appearance and geometry, by reducing the need to acquire environmental view-dependent effects. Our work leverages similar light-field characteristics to recover geometry more accurately within the volumetric rendering framework of NeRF, adding a prior to how geometry and appearance should present in the scene.

## 3. Method

Surf-NeRF introduces novel regularisation to locally enforce the properties of surface light fields. By regularising towards this representation, we represent smooth geometry with a physically viable separation of appearance and geometry producing more geometrically accurate radiance fields. This reduces shape radiance ambiguity by encouraging continuous regions of density with a smoothly varying view-dependent appearance. A visual depiction of our methodology is shown in Figure 2.

### 3.1. Preliminaries

A surface light field is a subset of the plenoptic function [9] that provides the colour of a light ray originating from a surface within the scene. Solid scene geometries are well represented given the decoupled parameterisation of the scene geometry and radiance on these surfaces.

A surface light field $L$ exists strictly on a surface geometry $G$ mapping directions in the unit 2-sphere ($\mathbf{S}^2$) to radiance, $L : G \times \mathbf{S}^2 \to \mathbf{c}$, for an RGB colour triplet $\mathbf{c}$ [41].

Few real surfaces have an entirely view-dependent appearance. Similar to traditional light fields [19], a surface light field may be decomposed into a Lambertian (diffuse) reflectance, $S$, and specular (or more generally a view-dependent) component, $R$ [24, 41], $L(\mathbf{u}, \boldsymbol{\omega}) = S(\mathbf{u}) + R(\mathbf{u}, \boldsymbol{\omega})$. These intrinsic decompositions are defined for a point on the surface, $\mathbf{u}$, and viewing direction, $\boldsymbol{\omega}$. The diffuse reflectance varies only with position over the surface, whilst the view-dependent component captures elements like reflection and refraction at the surface [20].

A surface light field has a strong assumption that the radiance seen from a given direction is provided entirely from a continuous geometry in the scene. This accordingly models phenomena such as volumetric scattering, reflection or transmission as a function of surface position rather than in a volumetric (NeRF) or physics-based (rendering engine) manner. In this way, a surface light field is decoupled from the geometry of a scene, meaning a surface may be deformed while its appearance looking along a ray is maintained [41]. Our approach is motivated by this decoupling to encourage Lambertian radiance to exist on smooth sheets of density, or surfaces, with a view-dependent colour.

### 3.2. Model

Our proposed method builds on the reflection parameterisation in Ref-NeRF [36], which splits the scene into a diffuse and specular appearance term similar to a surface light field. This parameterisation provides enhanced results across Lambertian and specular scenes over original NeRF variants [3, 27] by learning a spatially varying diffuse colour $\mathbf{c}_d$, specular tint $\mathbf{s}$, rendering normal $\hat{\mathbf{n}}'$ and a view-dependent specular colour $\mathbf{c}_s$ for each point in the scene. The final colour of a ray in this parameterisation is given as $\mathbf{c} = \mathbf{c}_d + \mathbf{s}\mathbf{c}_s$, a linear combination of the diffuse and specular terms. Ref-NeRF struggles to entirely separate Lambertian and specular appearance, utilising density instead to explain complex phenomena such as non-planar reflection, anisotropic reflection and interreflection, leading to the results seen in Figure 1. We apply our regularisation losses denoted by $\mathcal{L}$ to ZipNeRF [6], leveraging its state of the art performance with grid-based encoding [28], using the Ref-NeRF parameterisation. We make no modifications to the model itself beyond this, but include a second, data-driven sampling to impose physically-inspired regularisation. We maintain two proposal networks and one NeRF network with 64 and 32 samples per ray. Importantly, our regularisation terms are extensible to other NeRF variants, and may be applied after the main training as a fine-tuning
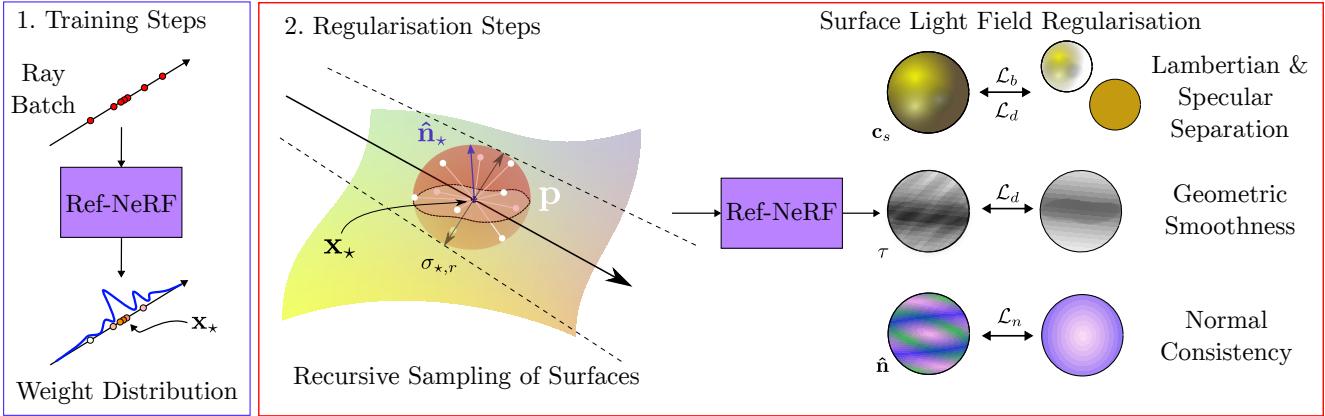
Figure 2. An overview of the Surf-NeRF methodology. We use a first surface assumption of light to locate samples $\mathbf{x}_\star$ which are likely to lie on a surface in the scene. Sampling at this point in multiple directions, and at points nearby using points drawn from a sphere $\mathbf{p}$, we impose regularisation on directional and spatial behaviour in this region. We separate the Lambertian component from the specular colour channel $\mathbf{c}_s$ by sampling points through multiple directions. Using neighbouring points we also regularise geometric smoothness on density $\tau$ and consistency of normals $\hat{\mathbf{n}}$ leading to improved, continuous geometry.

stage, as we show in Section 4.

### 3.3. Sampling Radiance at a Surface

We regularise the surface and its light field by sampling a batch of positions and directions at the point where the ray insects a surface in the scene, as shown in Figure 3. We formulate local regularisation terms based on surface light field properties and the current scene geometry and appearance at this location.

Where prior works have sampled unseen image patches [29] to encourage consistent depths, we sample a batch of unseen rays localised at a surface point in the scene to encourage local geometric continuity. This batch approach also has significant benefit over single perturbed points seen in prior work [30, 48], as it allows for changes in surface orientation and structure not captured by a single sample to be accounted for, as we detail in Section 3.4.

We utilise a first-surface assumption to infer the location of geometry; we assume the majority of radiance is emitted by dense points closest to the camera. Our candidate surface is the first point along a ray with weight $w(\mathbf{x}_\star)$ greater than the median weight of the ray $\text{med}(\mathbf{w}(\mathbf{x}))$, preventing regularisation from occurring behind the true location of the surface. This minimises sampling around points which are occluded and therefore not well positioned with multi-view constraints. Choosing the median ensures that this selection is more robust to skewed weight distributions, particularly early during training. This is shown in Figure 2 left.

Using the point $(\mathbf{x}_\star)$, origin $(\mathbf{o})$, direction $(\mathbf{d})$, surface normal $(\hat{\mathbf{n}}_\star)$, and covariance $(\boldsymbol{\Sigma}_\star)$ of this sample, we generate two new batches; a spatial batch to regularise density and a directional batch to regularise appearance.

We adapt a deterministic sampling scheme on the sphere [17] to produce uniformly distributed points used

in both regularisation batches. Samples drawn from a von Mises-Fischer distribution with concentration parameter $\kappa = 0$ are uniformly distributed on the unit 2-sphere $\mathbf{S}^2 \subset \mathbb{R}^3$. We sample this distribution deterministically using a Fibonacci-Kronecker lattice as proposed by [17] for $N$ samples. This is partitioned into $\log_2 N$ shells sampling the unit ball. Similar to the unscented sampling presented in ZipNeRF [6], we apply a random rotation [1] to these samples during training to avoid any bias from the orientation of the sampling scheme, arriving at samples $\mathbf{p}$. Further details are provided in the supplementary material.

To produce virtual rays which can be cast during training, we use our sphere samples to define directions $\mathbf{d}_s = \hat{\mathbf{p}}$, origins $\mathbf{o}_s$ and covariances $\boldsymbol{\Sigma}_s$ to construct new conical frusta. The origins and covariances for these rays,

$$\mathbf{o}_s = \mathbf{R}_s(\mathbf{o} - \mathbf{x}_\star) + \mathbf{x}_\star, \quad \boldsymbol{\Sigma}_s = \mathbf{R}_s \boldsymbol{\Sigma}_\star \mathbf{R}_s^\top, \quad \text{(1a,b)}$$

are defined by the rotation matrix $\mathbf{R}_s$ which rotates the original ray direction $\mathbf{d}$ to each $\mathbf{d}_s$. This results in $N$ rays at the same distance to the surface as the original ray, with MipN-eRF gaussians [5] aligned along ray directions.

We use these new virtual rays to produce two batches, querying the density field $\tau$ in the local 3D region of the surface at $\mathbf{x}_\star$ and the distribution of view-dependent colours $\mathbf{c}_s$ through a range of viewing angles. We refer to these as the *spatial* and *directional* batches. Below we outline the locations $\mathbf{x}$, and directions $\mathbf{d}$ through which we sample the NeRF in these batches. We also visually depict our sampling schemes in Figure 3.

**Directional Sampling:** Specularities at a surface are piecewise smooth, taking on the characteristics of what is being reflected. By sampling a point on a surface in a range of viewing directions and characterising the colour distribution, we can quantify how closely it matches the surface
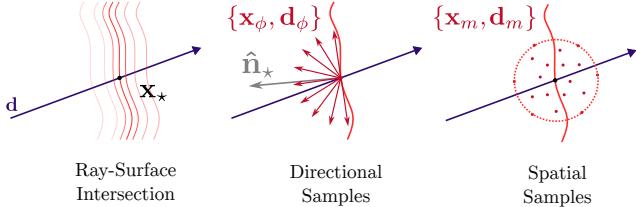
**Ray-Surface Intersection** — **Directional Samples** — **Spatial Samples**

Figure 3. A visual depiction of the two sampling batches localised at a surface in the scene. After defining a ray surface intersection $\mathbf{x}_\star$, we use uniform samples on a unit ball to sample the surface through multiple viewing angles, and the local 3D volume.

light field model introduced above. The directional batch $(\mathbf{x}_\phi, \mathbf{d}_\phi)$ samples at a single spatial location, but through a range of outward viewing angles:

$$\{\mathbf{x}_\phi,\ \mathbf{d}_\phi\} = \{\mathbf{x}_\star,\ \mathrm{sgn}(\mathbf{d}_s \cdot \mathbf{n}_\star)\mathbf{d}_s\}, \qquad (2)$$

where $\mathrm{sgn}(\cdot)$ is the sign function.

**Spatial Sampling:** In a NeRF, geometry is characterised by the density field $\tau(\mathbf{x})$, whose gradient with respect to position $\mathbf{x}$ approximates surface normals $\mathbf{n} \approx -\nabla_\mathbf{x}\tau$. By sampling density and normals of neighbouring points, and enforcing consistency between these values we can provide additional geometric supervision towards a surface model. The spatial batch $(\mathbf{x}_m, \mathbf{d}_m)$ consists of points $\mathbf{p}$ in the unit ball located at the surface $\mathbf{x}_\star$,

$$(\mathbf{x}_m,\ \mathbf{d}_m) = \left(\mathbf{x}_\star + \sigma_{\star,r}\mathbf{p}\sqrt{2\ln 2},\ \mathbf{d}_s\right), \qquad (3)$$

where $\sigma_{\star,r}$ is the radial variance of the sample as per mip-NeRF [3]. This samples the NeRF within a single pixel's conical frustum, ensuring regularisation occurs at the scale of the training data; images closer to the scene regularise at finer detail compared to those further away. Importantly, our sample volume adapts in scale towards a minimum as the NeRF localises density during training during the later stages of training. We characterise this behaviour in the supplement. As we do not care about the colour of these samples off the surface, the spatial batch uses the directions $\mathbf{d}_s$ of the virtual rays.

### 3.4. Local Smoothness

We regularise the geometry of the representation by penalising sparse and irregular density, allowing for a smooth and continuous sheet-like density field to form. These constraints are valid except where topological edges or corners exist; we therefore formulate $L_1$-norm regularisation terms to allow for these local features to form when required.

The spatial sampling batch encompasses the region in front of and behind the candidate surface. Using the normal at the candidate surface, we penalise points proportional to their density $\tau_j$ and perpendicular distance $\left|\frac{(\mathbf{x}_{m,j} - \mathbf{x}_\star)}{\|(\mathbf{x}_{m,j} - \mathbf{x}_\star)\|} \cdot \hat{\mathbf{n}}_\star\right|$ from the surface encouraging a plane of

density to form perpendicular to the normal. Points far from the surface should have low density, whilst those on the same plane as the sample point should be dense. This loss term,

$$\mathcal{L}_d = \lambda_d \sum_{\mathbf{r} \in \mathbf{R}} w_\star \sum_{j=0}^{N-1} \left(1 - e^{-\tau_j}\right)\left|\frac{(\mathbf{x}_{m,j} - \mathbf{x}_\star)}{\|(\mathbf{x}_{m,j} - \mathbf{x}_\star)\|} \cdot \hat{\mathbf{n}}_\star\right|, \tag{4}$$

sums over $N$ samples in the spatial batch and is also weighted by the surface point $w_\star$ rendering weight. This ensures that early during training, the NeRF is able to change surface locations as more rays sample the scene. We do not place a stop gradient on $\hat{\mathbf{n}}_\star$, allowing for the surface to re-orient based on the density of samples in the local volume.

Normals in this local region should exhibit similar smoothness, and so the term,

$$\mathcal{L}_n = \lambda_n \sum_{\mathbf{r} \in \mathbf{R}} w_\star \sum_{j=0}^{N-1} \left(1 - e^{-\tau_j}\right)\frac{(1 - (\hat{\mathbf{n}}_j \cdot \hat{\mathbf{n}}_\star))}{2}, \tag{5}$$

is imposed proportional to the angle between the normal at $\mathbf{x}_\star$ and those at each sample point $\mathbf{p}_j$. This term encourages dense samples to have a normal which is parallel with that at the surface point. Both terms localise density and encourage smooth, continuous geometries. We weight these terms by hyperparameters $\lambda_d$ and $\lambda_n$ respectively.

### 3.5. Visually Realistic Appearance Separation

Whilst highly specular surfaces or highly Lambertian surfaces are well represented with a reflection encoding, surfaces with both a strong Lambertian and specular component exhibit shape-radiance ambiguity. This is characterised by incorrect geometry with a view-dependent colour bias; that is the minimum specular colour over all viewing angles is non-zero to satisfy the appearance in all training view-points. More details are found in the supplementary.

To encourage the surface light field separation of the Lambertian and view-dependent terms we propose two additional losses based on directionally sampling a point. To encourage a minimum specular bias, we penalise the normalised specular colour over viewing angles in the directional sampling batch $\mathbf{c}_{s,\phi}$,

$$\mathcal{L}_b = \lambda_b \sum_{\mathbf{r} \in \mathbf{R}} w_\star \sum_{j=0}^{N-1} \left\|\frac{\mathbf{c}_{s,j}}{\mathrm{sg}(\|\max_\phi \mathbf{c}_s\|)}\right\|^2, \tag{6}$$

where $\mathrm{sg}(\cdot)$ is the stop gradient operator. Where the specular distribution is uniform over viewing angles, this loss term is maximal. This ensures that the appearance of a point must accumulate most of its signal in either the specular or the Lambertian term - not in both. To maintain a piecewise

smooth distribution in the specular, we enforce a total variation loss over the specular colour:

$$\mathcal{L}_s = \lambda_s \sum_{\mathbf{r} \in \mathbf{R}} w_\star \sum_{j=0}^{N-1} \mathrm{STV}_j(\mathbf{d}_\phi, \mathbf{c}_s). \tag{7}$$

This total variation is applied over the surface of the sampling sphere as a form of graph total variation. This term,

$$\mathrm{STV}_j(\mathbf{d}_\phi, \mathbf{c}_s) = \sum_k \frac{1}{2}(\mathbf{d}_j \cdot \mathbf{d}_k + 1)\|\mathbf{c}_j - \mathbf{c}_k\|_1, \tag{8}$$

is based upon the $k$-Nearest Neighbours of each sample point and weighted by the cosine distance between the sample and its neighbours. This approach provides a localised consistency which is edge-aware continuous over the sampling sphere, capturing occlusions. Both of these terms smooth the specular distribution and increase the reliance on the diffuse component of the model and are weighted by hyperparameters $\lambda_b$ and $\lambda_s$ respectively.

### 3.6. Curriculum Learning

We impose our surface regularisation under curriculum learning to maintain visual fidelity whilst still improving geometry. As training progresses and additional regularisation is imposed more frequently, the representation is refined to satisfy the geometric and appearance characteristics of a surface light field. The surface constraints presented in Sections 3.4-3.5 are of similar computational complexity to a regular training step, requiring a second pass through both MLPs and therefore taking close to double the time. By using curriculum learning we trade-off the improvement to the representation with training time.

We schedule regularisation to occur in a staircase schedule of powers of 2: every 512 iterations early in the training, down to every 4 iterations in the final stages. On this schedule, there is approximately a 25% increase to training time.

For each regularisation sample, we add on each regularisation loss $\mathcal{L}$ to the losses for the combined ZipNeRF and Ref-NeRF model structure. More model implementation details are provided in the supplementary material, along with hyperparameter weights.

## 4. Results

Our regularisation is implemented in JAX [13] and based on the ZipNeRF codebase [7]. We utilise the improvements to quality and speed in ZipNeRF combined with the physics-based model structure introduced in Ref-NeRF. We compare to prior works MipNeRF360 [4] and Ref-NeRF [36]. We evaluate Surf-NeRF on the Shiny Real and the Shiny Objects [36] datasets used in Ref-NeRF as the main benchmark for our approach. We also evaluate on

a captured dataset consisting of 4 complex reflective objects under controlled illumination, called the Koala dataset. Each scene contains approximately 40 training images and 10 test images with ground truth poses obtained using a Universal Robots UR5e robotic arm. We select objects that include non-planar specularities, fine details and specularities which present as virtual images in front of the surface. Further details are provided in the supplementary material.

### 4.1. Shiny Objects Dataset

Table 1 summarises results on the Shiny Objects dataset, comparing both visual fidelity by peak-signal-to-noise ratio (PSNR) in decibels (dB) and structural similarity scores (SSIM), and geometric fidelity by mean angular error (MAE) of the density derived normals in degrees and root mean squared error (RMSE) of the disparity in inverse units. We compare across positionally encoded baselines based on MipNeRF360 [4], denoted as "PE-based" and "grid-based" implementations based on ZipNeRF [6]. We also compare against a variant of MipNeRF which includes a diffuse colour channel, denoted as "MipNeRF+Diff". We train the positionally encoded implementations using V2-8 TPUs, and the grid-based methods on four NVIDIA V100 32GB GPUs. Per scene results are included in the supplementary materials.

Table 1. Summary on the *Shiny Objects* [36] dataset between visual metrics (PSNR, SSIM) and geometric metrics (MAE, RMSE). Our method yields more accurate normals, trading off a maginal decrease to PSNR. Yellow, orange and red are the third, second, and first scores.

| | Model | PSNR ↑ | SSIM ↑ | MAE ↓ | RMSE ↓ |
|---|---|---|---|---|---|
| PE | MipNeRF | 31.29 | 0.943 | 56.17 | 0.125 |
| | MipNeRF+Diff | 30.84 | 0.936 | 60.00 | 0.122 |
| | Ref-NeRF | 33.21 | 0.971 | 25.89 | 0.117 |
| | Surf-NeRF (Ours) | 33.01 | 0.967 | 22.15 | 0.117 |
| Grid | ZipNeRF | 31.02 | 0.952 | 19.47 | 0.209 |
| | Zip+Ref-NeRF | 33.14 | 0.964 | 14.70 | 0.195 |
| | Surf-NeRF (Ours) | 32.13 | 0.954 | 13.35 | 0.182 |

Across both the positional encoding based and hash based networks we see comparable performance to Ref-NeRF in visual fidelity, demonstrating our ability to maintain visual fidelity whilst encouraging a more geometrically accurate representation. In the positionally encoded approaches we see a 14.4% improvement to normal accuracy, with notable improvements around regions of high specularity. We also see a 9.2% improvement to normals for the grid based implementation. Given the discretisation of the grid-based approaches, we see comparatively higher disparity errors across the dataset but improved normals from the interpolation of encodings. Our approach realises a 6.7% improvement in disparity versus Zip+Ref-NeRF.

In Figure 4 we demonstrate that the proposed regularisation is able to correct for warping of geometry and improve its consistency on the car scene, whilst also cor-

rectly separating the Lambertian racing stripes. Similarly, on the toaster dataset we successfully separate the Lambertian toast whilst retaining its reflection and improving geometry around specular materials with high curvature.
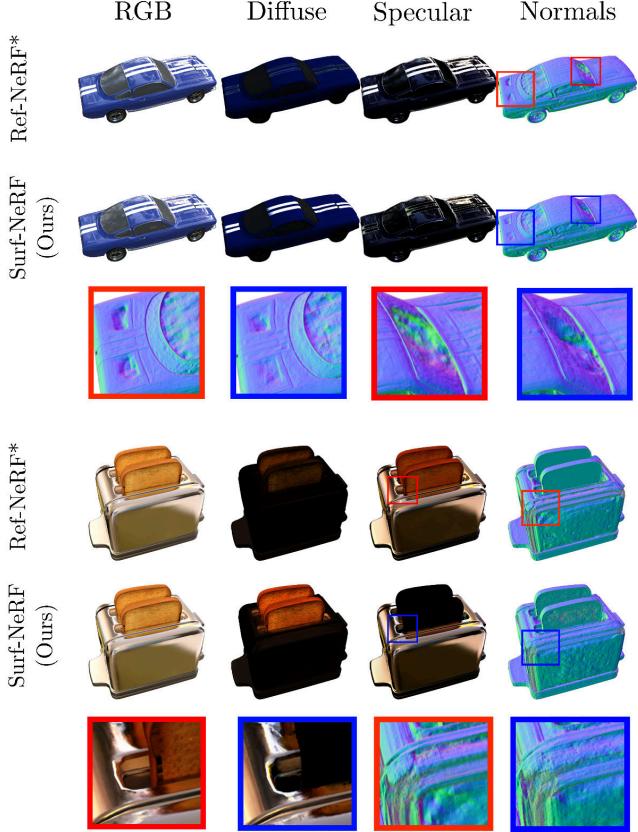


Figure 4. Results from the grid based models. By regularising the scene, we not only remove Lambertian scene content from the specular, like the stripes on the car and toast in the toaster, but our approach yields improved geometry particularly around curves and specularity, like the car windshield and corners of the toaster. Ref-NeRF* indicates the Zip+Ref-NeRF model.

Figure 5 depicts performance of positionally encoded models on specular surfaces with low texture. Surf-NeRF yields drastically better surface geometry and normals.

## 4.2. Shiny Real Dataset

The visual fidelity results on the *Shiny Real* dataset [36] are shown in Table 2 for Surf-NeRF and comparison baseline models. We demonstrate comparable visual fidelity to other state of the art NeRF variants demonstrating no appreciable drop in visual fidelity compared to non-regularised approaches on the positional encoding based models. The complex appearances of these datasets combined with the hash-encoding sees a small drop in PSNR values for grid-based variants. We show decreased reliance on floaters to explain appearance via Surf-NeRF in the supplement.
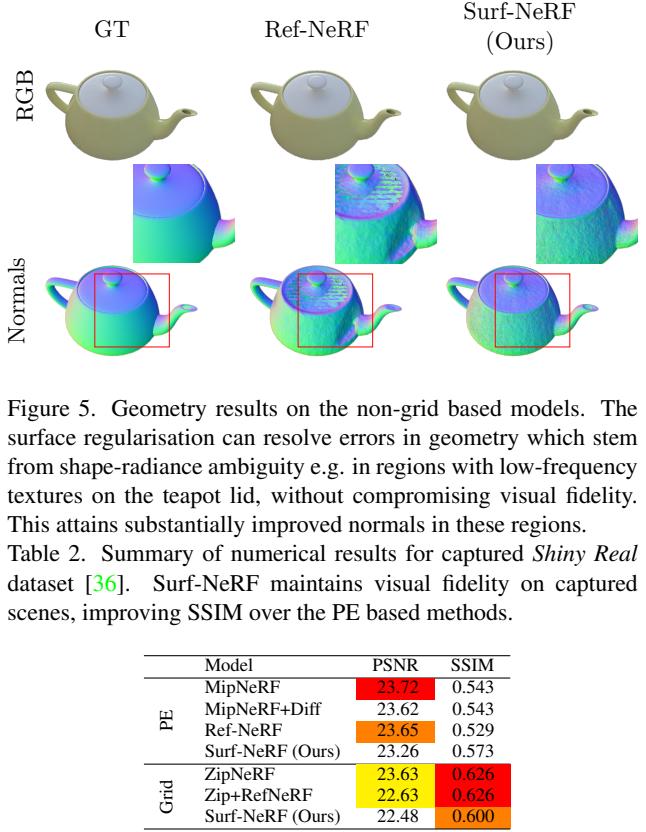


Figure 5. Geometry results on the non-grid based models. The surface regularisation can resolve errors in geometry which stem from shape-radiance ambiguity e.g. in regions with low-frequency textures on the teapot lid, without compromising visual fidelity. This attains substantially improved normals in these regions.

Table 2. Summary of numerical results for captured *Shiny Real* dataset [36]. Surf-NeRF maintains visual fidelity on captured scenes, improving SSIM over the PE based methods.

| | Model | PSNR | SSIM |
|---|---|---|---|
| PE | MipNeRF | 23.72 | 0.543 |
| | MipNeRF+Diff | 23.62 | 0.543 |
| | Ref-NeRF | 23.65 | 0.529 |
| | Surf-NeRF (Ours) | 23.26 | 0.573 |
| Grid | ZipNeRF | 23.63 | 0.626 |
| | Zip+RefNeRF | 22.63 | 0.626 |
| | Surf-NeRF (Ours) | 22.48 | 0.600 |

## 4.3. Koala Dataset

Table 3 shows the strength of our proposed regularisation on complex reflective scenes. Figure 6 demonstrates the ability for the proposed regularisation to separate Lambertian and specular colours in a more physically consistent manner, and to overcome failure cases where specularities are not densely sampled by the training data. Additional results are provided in the supplementary.

Table 3. Summary of numerical results for captured *Koala* dataset. Surface regularisation helps scene convergence around complex specular geometry improving PSNR.

| Model | PSNR | SSIM |
|---|---|---|
| ZipNeRF | 27.79 | 0.663 |
| Zip+RefNeRF | 23.69 | 0.632 |
| Surf-NeRF (Ours) | **27.24** | **0.655** |

## 4.4. Ablation Study

In Table 4 we present an ablation study of our additional surface regularisation terms on the Helmet scene from the Shiny Objects dataset. As shown, including our full approach provides strong performance across both the visual and geometric metrics. Our appearance based regularisation is coupled closely to the geometry of the scene. Removing the specular bias term drastically increases the normal error, however does not alter disparity in the scene. The normal term improves continuity of the geometry in the scene.
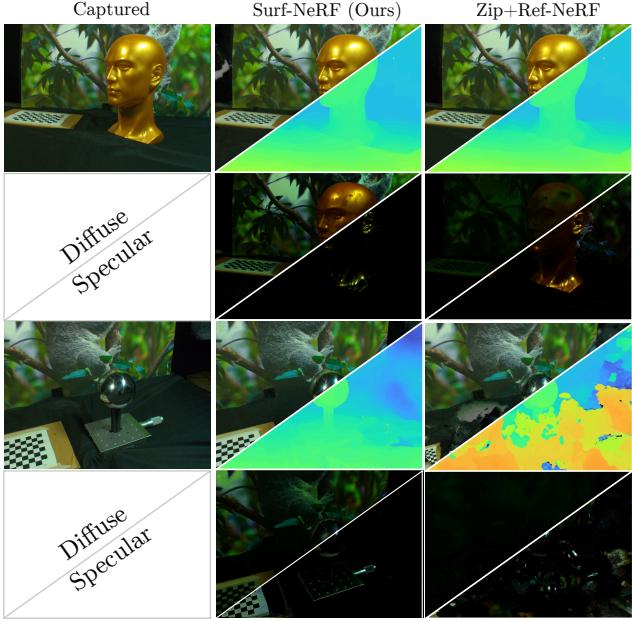
Figure 6. The "gold head" and "shiny ball" scenes from our captured *Koala* dataset. The first row contains RGB and median depth renderings, whilst the second contains the diffuse and specular renderings. The Gold Head dataset shows complete separation of Lambertian and specular appearance for the mannequin. The Shiny Ball dataset represents a failure case for the Zip+Ref-NeRF baseline, however the proposed regularisation regularises these floaters away during training and appropriately places the majority of the colour into the diffuse colour term.

Whilst removing it allows for each individual normal to be determined more freely, reducing the MAE, it also allows for floaters to be introduced unnecessarily by our spatial term corresponding to an increase in PSNR but also an increase in disparity error within the scene.

Table 4. Ablation study on Surf-NeRF regularisation terms on the Helmet scene. All terms contribute to overall performance.

| | $\mathcal{L}_s$ | $\mathcal{L}_n$ | $\mathcal{L}_d$ | $\mathcal{L}_b$ | PSNR↑ | SSIM↑ | MAE↓ | RMSE↓ |
|---|---|---|---|---|---|---|---|---|
| (Ours) | ✓ | ✓ | ✓ | ✓ | 32.93 | 0.974 | 8.906 | 0.192 |
| w/o spatial | | ✓ | ✓ | ✓ | 32.33 | 0.972 | 10.67 | 0.196 |
| w/o normal | ✓ | | ✓ | ✓ | 32.86 | 0.976 | 7.98 | 0.221 |
| w/o spec. TV | ✓ | ✓ | | ✓ | 32.15 | 0.973 | 10.01 | 0.195 |
| w/o bias | ✓ | ✓ | ✓ | | 32.38 | 0.973 | 16.05 | 0.192 |

Table 5 demonstrates an study on our curriculum learning approach. We demonstrate the effectiveness of our approach as a function of the frequency of regularisation. We

Table 5. Parameter study on regularisation frequency on Coffee scene. There is a tradeoff to regularisation frequency between quality, time and effect, as in Figure 1.

| Initial Freq. | Final Freq. | PSNR ↑ | SSIM ↑ | MAE ↓ | RMSE ↓ |
|---|---|---|---|---|---|
| 1024 | 8 | 31.38 | 0.967 | 14.89 | 0.181 |
| 512 | 4 | 31.99 | 0.968 | 15.29 | 0.181 |
| 256 | 2 | 31.70 | 0.964 | 18.17 | 0.208 |
| 128 | 1 | 32.09 | 0.967 | 15.95 | 0.181 |

demonstrate improvement to the normals and disparity in all cases compared to Zip+Ref-NeRF. Note that performance is highest for regularisation at a rate which is sparse enough for adaptation to changing density and frequent enough that corrections to the representation are retained.

## 4.5. Surface Regularisation as Finetuning

We demonstrate surface regularisation as a fine-tuning step, after a NeRF has been trained on pre-trained Zip+Ref-NeRF and ZipNeRF model on the *car* scene from the Shiny Objects dataset. As ZipNeRF has no Lambertian colour term, we do not include $\mathcal{L}_b$, and apply our sphere total variation regularisation $\mathcal{L}_s$ to the view-dependent colour **c**, an example of our applicability to models without reflection parameterisation. The spatial sampling terms depend only on the geometry of the NeRF and so may be applied to other density-based neural fields without modification. We continue training at a fixed learning rate of $3.25 \times 10^{-4}$ and present our results in Table 6. In both cases we demonstrate an improvement in normals over both the base model and the naïve approach where the model continues training without our regularisation. Note that finetuning with the reflection parameterisation, as in our Zip+RefNeRF model, results in higher disparity errors as the NeRF changes density placement during the additional training.

Table 6. Surf-NeRF as a finetuning step and applicability to other NeRF variants. With little additional training, we refine surface normals.

| Base | Extra | PSNR | SSIM | Median MAE | RMSE |
|---|---|---|---|---|---|
| | - | **27.34** | **0.932** | 24.21 | 0.220 |
| ZipNeRF | 2.5k ZipNeRF | 27.32 | 0.931 | 24.30 | 0.220 |
| | 2.5k SurfNeRF | 27.33 | 0.931 | **23.57** | **0.218** |
| | - | **30.25** | **0.957** | 18.44 | **0.246** |
| ZipRefNeRF | 2.5k ZipRefNeRF | 30.14 | 0.956 | 18.59 | 0.251 |
| | 2.5k SurfNeRF | 30.13 | 0.956 | **18.29** | 0.248 |

## 4.6. Limitations

Modelling the scene as surfaces, elements that are better represented volumetrically like hair or subsurface scattering may be unnecessarily regularised and their representation degraded using our proposed regularisation. Since Surf-NeRF makes use of a second pass through the network, our methodology takes twice as long on each regularisation step compared to a vanilla training step. We enforce geometric consistency during training at the cost of restricting the networks ability to deliver the best photometric accuracy, leading to marginally reduced PSNR scores.

## 5. Conclusion

We have proposed Surf-NeRF, a novel regularisation approach and four novel physically derived regularisers, and shown these to improve geometry in a NeRF by curricu-

lum learning towards a surface light field model. This ensures consistency of density and normals on geometry in the volume, in addition to a physics-inspired regularisation on view-dependent appearance. This enables a NeRF to adjust geometric structure during training without greatly sacrificing visual fidelity, demonstrating qualitatively improved continuity of geometry, with improvements of up to 14.4% in normals and 6.7% in disparity.

This work produces better NeRF representations for geometry critical tasks with unchanged model architectures.

# References

[1] James Arvo. Fast random rotation matrices. In David Blair Kirk, editor, *Graphics Gems III (IBM Version)*, The Graphics Gems Series, pages 117–120. Academic Press, 1992. 4

[2] Dejan Azinović, Ricardo Martin-Brualla, Dan B Goldman, Matthias Nießner, and Justus Thies. Neural RGB-D Surface Reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6290–6301, 2022. 2

[3] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-NeRF: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021. 2, 3, 5

[4] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-NeRF 360: Unbounded Anti-Aliased Neural Radiance Fields. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 6

[5] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-NeRF 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022. 4

[6] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Zip-NeRF: Anti-Aliased Grid-Based Neural Radiance Fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 19697–19705, October 2023. 2, 3, 4, 6

[7] Jonathan T. Barron, Keunhong Park, Ben Mildenhall, John Flynn, Dor Verbin, Pratul Srinivasan, Peter Hedman, Philipp Henzler, and Ricardo Martin-Brualla. CamP Zip-NeRF: A Code Release for CamP and Zip-NeRF, 2024. 6

[8] Mojtaba Bemana, Karol Myszkowski, Jeppe Revall Frisvad, Hans-Peter Seidel, and Tobias Ritschel. Eikonal Fields for Refractive Novel-view Synthesis. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–9, 2022. 2

[9] James R Bergen and Edward H Adelson. The Plenoptic Function and The Elements of Early Vision. *Computational Models of Visual Processing*, 1:8, 1991. 3

[10] Mark Boss, Raphael Braun, Varun Jampani, Jonathan T Barron, Ce Liu, and Hendrik Lensch. NeRD: Neural Reflectance Decomposition from Image Collections. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12684–12694, 2021. 3

[11] Mark Boss, Andreas Engelhardt, Abhishek Kar, Yuanzhen Li, Deqing Sun, Jonathan T. Barron, Hendrik Lensch, and Varun Jampani. SAMURAI: Shape And Material from Unconstrained Real-world Arbitrary Image collections. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. 3

[12] Mark Boss, Varun Jampani, Raphael Braun, Ce Liu, Jonathan Barron, and Hendrik Lensch. Neural-PIL: Neural Pre-Integrated Lighting for Reflectance Decomposition. *Advances in Neural Information Processing Systems*, 34:10691–10704, 2021. 3

[13] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. 6

[14] Anpei Chen, Zexiang Xu, Xinyue Wei, Siyu Tang, Hao Su, and Andreas Geiger. Factor Fields: A Unified Framework for Neural Fields and Beyond. *arXiv preprint arXiv:2302.01226*, 2023. 2

[15] Ziang Cheng, Hongdong Li, Richard Hartley, Yinqiang Zheng, and Imari Sato. Diffeomorphic Neural Surface Parameterization for 3D and Reflectance Acquisition. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–10, 2022. 3

[16] Andreea Dogaru, Andrei-Timotei Ardelean, Savva Ignatyev, Egor Zakharov, and Evgeny Burnaev. Sphere-guided training of neural implicit surfaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20844–20853, June 2023. 2

[17] Daniel Frisch and Uwe D. Hanebeck. Deterministic Von Mises–Fisher sampling on the sphere using fibonacci lattices. In *2023 IEEE Symposium Sensor Data Fusion and International Conference on Multisensor Fusion and Integration (SDF-MFI)*, pages 1–8, 2023. 4

[18] Qiancheng Fu, Qingshan Xu, Yew-Soon Ong, and Wenbing Tao. Geo-Neus: Geometry-Consistent Neural Implicit Surfaces Learning for Multi-view Reconstruction. In *Advances in Neural Information Processing Systems*, 2022. 2

[19] Elena Garces, Jose I Echevarria, Wen Zhang, Hongzhi Wu, Kun Zhou, and Diego Gutierrez. Intrinsic Light Field Images. In *Computer Graphics Forum*, volume 36, pages 589–599. Wiley Online Library, 2017. 3

[20] Elena Garces, Carlos Rodriguez-Pardo, Dan Casas, and Jorge Lopez-Moreno. A Survey on Intrinsic Images: Delving Deep into Lambert and Beyond. *Int. J. Comput. Vision*, 130(3):836–868, mar 2022. 3

[21] Wenhang Ge, Tao Hu, Haoyu Zhao, Shu Liu, and Ying-Cong Chen. Ref-NeUS: Ambiguity-reduced neural implicit surface learning for multi-view reconstruction with reflection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4251–4260, 2023. 2

[22] Yuan-Chen Guo, Di Kang, Linchao Bao, Yu He, and Song-Hai Zhang. NeRFReN: Neural Radiance Fields with Reflec-

tions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18409–18418, 2022. 2

[23] Mijeong Kim, Seonguk Seo, and Bohyung Han. InfoNeRF: Ray entropy minimization for few-shot neural volume rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12912–12921, 2022. 2

[24] Zhengqi Li and Noah Snavely. CGIntrinsics: Better Intrinsic Image Decomposition through Physically-Based Rendering. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018. 3

[25] Li Ma, Vasu Agrawal, Haithem Turki, Changil Kim, Chen Gao, Pedro Sander, Michael Zollhöfer, and Christian Richardt. SpecNeRF: Gaussian directional encoding for specular reflections. *arXiv preprint arXiv:2312.13102*, 2023. 2

[26] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. NeRF in the Wild: Neural Radiance Fields for Unconstrained Photo Collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7210–7219, 2021. 2

[27] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In *ECCV*, 2020. 1, 2, 3

[28] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, July 2022. 3

[29] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. RegNeRF: Regularizing Neural Radiance Fields for View Synthesis From Sparse Inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5480–5490, 2022. 2, 3, 4

[30] Michael Oechsle, Songyou Peng, and Andreas 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. 2, 4

[31] Daniel Rebain, Mark Matthews, Kwang Moo Yi, Dmitry Lagun, and Andrea Tagliasacchi. LolNeRF: Learn From One Look. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1558–1567, 2022. 2

[32] Sara Sabour, Suhani Vora, Daniel Duckworth, Ivan Krasin, David J. Fleet, and Andrea Tagliasacchi. Robustnerf: Ignoring distractors with robust losses. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20626–20636, 2023. 2

[33] Nagabhushan Somraj, Adithyan Karanayil, and Rajiv Soundararajan. SimpleNeRF: Regularizing sparse input neural radiance fields with simpler solutions. In *SIGGRAPH Asia 2023 Conference Papers*, SA '23, New York, NY, USA, 2023. Association for Computing Machinery. 3

[34] Matthew Tancik, Vincent Casser, Xinchen Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. Block-NeRF: Scalable Large Scene Neural View Synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8248–8258, 2022. 2

[35] A. Tewari, J. Thies, B. Mildenhall, P. Srinivasan, E. Tretschk, W. Yifan, C. Lassner, V. Sitzmann, R. Martin-Brualla, S. Lombardi, T. Simon, C. Theobalt, M. Nießner, J. T. Barron, G. Wetzstein, M. Zollhöfer, and V. Golyanik. Advances in Neural Rendering. *Computer Graphics Forum (EG STAR 2022)*, 2022. 2

[36] Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T Barron, and Pratul P Srinivasan. Ref-NeRF: Structured View-Dependent Appearance For Neural Radiance Fields. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5481–5490. IEEE, 2022. 1, 2, 3, 6, 7

[37] Dor Verbin, Pratul P Srinivasan, Peter Hedman, Ben Mildenhall, Benjamin Attal, Richard Szeliski, and Jonathan T Barron. NeRF-Casting: Improved view-dependent appearance with consistent reflections. *arXiv preprint arXiv:2405.14871*, 2024. 2

[38] Delio Vicini, Sébastien Speierer, and Wenzel Jakob. Differentiable Signed Distance Function Rendering. *ACM Transactions on Graphics (TOG)*, 41(4):1–18, 2022. 2

[39] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. NeuS: Learning Neural Implicit Surfaces by Volume Rendering for Multi-view Reconstruction. *Advances in Neural Information Processing Systems*, 34:27171–27183, 2021. 2

[40] Yiming Wang, Qin Han, Marc Habermann, Kostas Daniilidis, Christian Theobalt, and Lingjie Liu. NeuS2: Fast Learning of Neural Implicit Surfaces for Multi-view Reconstruction. *arXiv preprint arXiv:2212.05231*, 2022. 2

[41] Daniel N. Wood, Daniel I. Azuma, Ken Aldinger, Brian Curless, Tom Duchamp, David H. Salesin, and Werner Stuetzle. Surface Light Fields for 3D Photography. In *Proceedings of the 27th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '00, page 287–296, USA, 2000. ACM Press/Addison-Wesley Publishing Co. 2, 3

[42] Bangbang Yang, Chong Bao, Junyi Zeng, Hujun Bao, Yinda Zhang, Zhaopeng Cui, and Guofeng Zhang. Neumesh: Learning Disentangled Neural Mesh-Based Implicit Field for Geometry and Texture Editing. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XVI*, pages 597–614. Springer, 2022. 3

[43] Yao Yao, Jingyang Zhang, Jingbo Liu, Yihang Qu, Tian Fang, David McKinnon, Yanghai Tsin, and Long Quan. NeILF: Neural Incident Light Field for Physically-Based Material Estimation. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXI*, pages 700–716. Springer, 2022. 3

[44] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume Rendering of Neural Implicit Surfaces. *Advances in Neural Information Processing Systems*, 34:4805–4815, 2021. 2

10

[45] Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. MonoSDF: Exploring Monocular Geometric Cues for Neural Implicit Surface Reconstruction. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. 2

[46] Jason Zhang, Gengshan Yang, Shubham Tulsiani, and Deva Ramanan. NeRS: Neural Reflectance Surfaces for Sparse-view 3D Reconstruction in the Wild. *Advances in Neural Information Processing Systems*, 34:29835–29847, 2021. 3

[47] Kai Zhang, Fujun Luan, Zhengqi Li, and Noah Snavely. IRON: Inverse Rendering by Optimizing Neural SDFs and Materials from Photometric Images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5565–5574, 2022. 3

[48] Xiuming Zhang, Pratul P Srinivasan, Boyang Deng, Paul Debevec, William T Freeman, and Jonathan T Barron. Nerfactor: Neural Factorization of Shape and Reflectance Under an Unknown Illumination. *ACM Transactions on Graphics (TOG)*, 40(6):1–18, 2021. 4

[49] Bingfan Zhu, Yanchao Yang, Xulong Wang, Youyi Zheng, and Leonidas Guibas. VDN-NeRF: Resolving shape-radiance ambiguity via view-dependence normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 35–45, June 2023. 3

# Surf-NeRF: Surface Regularised Neural Radiance Fields - Supplementary Materials

## A. Deterministic Uniform Sphere Sampling

In this section we provide additional details regarding our uniform sphere sampling and its implementation. Sampling over spherical domains requires careful construction to avoid over-sampling or under-sampling, particularly as the azimuth coordinate $\theta \to \pm\frac{\pi}{2}$. We experimented extensively, and found that traditional uniform sampling leads to unintended clusters of sampling directions which led to artefacts in the specular term. For this reason, we adapt the approach proposed by Frisch et al. [17] designed to sample a von Mises-Fisher (vMF) distribution in $\mathbb{S}^2$ to provide a deterministic uniform sample of a sphere.

The vMF distribution,

$$f(\mathbf{x}) = \frac{\kappa}{2\pi(e^{\kappa} - e^{-\kappa})} e^{\kappa \boldsymbol{\mu}^{\top} \mathbf{x}}, \qquad (1)$$

is a probability distribution with mean $\boldsymbol{\mu} \in \mathbb{S}^2$ and concentration parameter $\kappa \geq 0$ on the 2-sphere. In the limit as $\kappa = 0$, $f(\mathbf{x}) = (4\pi)^{-1}$, resulting in a uniform distribution.

Frisch et al. [17] sample this distribution using a Fibonacci-Kronecker lattice, with coordinates $\mathbf{x}_{FK} \in \mathbb{S}^2$,

$$\mathbf{x}_{F,i} = \begin{bmatrix} w \\ \sqrt{1-w^2} \cos \frac{2\pi i}{\Phi} \\ \sqrt{1-w^2} \sin \frac{2\pi i}{\Phi} \end{bmatrix}, \ i \in \{0, \ldots, N-1\} \quad (2)$$

where $\frac{1}{\Phi} = \frac{\sqrt{5}-1}{2}$ is the Golden Ratio, and

$$w = 1 + \frac{1}{\kappa} \log 1\mathrm{p}\left(\frac{2i-1}{2N} \mathrm{expm1}(-2\kappa)\right), \qquad (3)$$

for $\log 1\mathrm{p}(x) = \log(1+x)$ and $\mathrm{expm1}(x) = \exp(x) - 1$. Since we want a uniform distribution, we take the limit $\kappa \to 0$, resulting in,

$$\lim_{\kappa \to 0} w = w_{\kappa=0} = \frac{1 - 2i + N}{N}. \qquad (4)$$

We wish to maintain the distribution of directions of the sphere, but partition the samples through the volume of the unit ball. We introduce a discrete radius for each sample,

$$\mathbf{x}_i = \frac{1}{\log_2 N}(1 + i \bmod \log_2 N)\mathbf{x}_{F,i}, \qquad (5)$$

to partition the samples into $\log_2 N$ shells in the unit ball.

We recognise that these shells, whilst uniform in direction have decreasing density with distance to the centre of the sphere. This helps in our use case for the sampling, with a higher density of samples closer to the centre of the ball (i.e. the surface). To avoid any bias which may arise in the volumetric samples, we generate a rotation matrix which is uniformly distributed in $SO(3)$, using [1]. Figure A.1 provides a visualisation of the direction and spatial coordinates for each sample using this approach.

We initialise a sphere of $N$ (in our case, $N = 32$) samples prior to the training loop, and at each iteration generate a unique rotation matrix for each ray which is applied to this sphere to provide the unique sampling. Rays which view the same part of the scene are therefore unlikely to sample to same positions or directions during training.
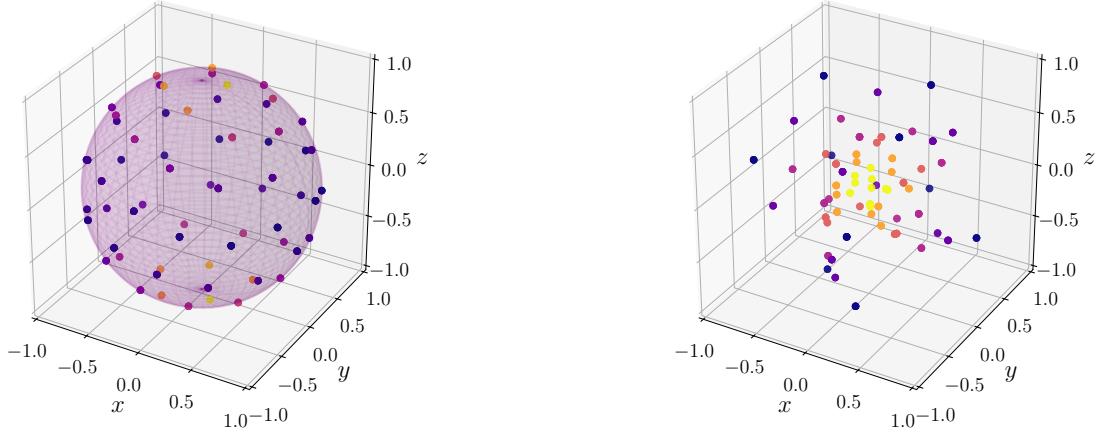
In Table A.1 we present a short study comparing random uniform sampling of points on the sphere to our proposed deterministic sampling. Our deterministic sampling scheme achieves approximately a 13.2% improvement to normals compared to a 5% improvement attained by uniform samples on the sphere. The degradation in PSNR and disparity can be attributed to clustering of sample points leading to large regions which are sparsely regularised during training, and regions which are more heavily regularised.

Table A.1. Deterministic sampling scheme study on the car scene from *Shiny Objects* [36] dataset. Deterministic samples improve in all aspects by preventing clustering with our comparatively small number of additional samples in a regularisation batch.

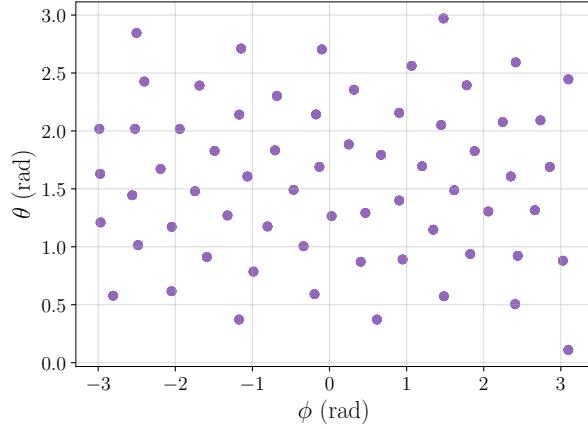| Method | PSNR | SSIM | MAE | Disp. RMSE |
|---|---|---|---|---|
| Zip+Ref-NeRF | 30.30 | 0.957 | 15.167 | 0.209 |
| Ours, Uniform Random | 29.67 | 0.953 | 14.668 | 0.216 |
| Ours, Deterministic | **29.94** | **0.954** | **13.165** | **0.198** |

## B. Lambertian Bias

Whilst positional encoding helps MLPs to learn higher frequency information [27] and reflection encodings [25, 36, 37] encourage a separation of low-frequency colour information with view-dependent components, we find in practice that the view-dependent term maintains a view-independent component across viewing angles. This is visu-

(a) 3D plot of deterministic samples on the sphere before applying the $\log_2 N$ radii partitions to sample the unit ball.

(b) 3D plot of deterministic sampling. These point locations are used in our proposed spatial sampling batch.



(c) Directions through which we sample. We see that the lattice uniformly distributes sample directions, taking into account the warping close to the poles.

Figure A.1. Example directional batch and its polar coordinate map, indicating samples used for regularisation.

alised in Figure B.1. Early in training this phenomenon can be useful, allowing for density to accumulate independent of the diffuse colour, and for surfaces with high roughness to develop low frequency specular distributions. In many cases, these low frequency components remain in the view-dependent appearance of the scene given minimal photometric error and no loss regularising this component. This has the drawback however that the representation does not need to assign a constant view-independent colour to points in the scene allowing for ambiguous placement of scene geometry. For most use cases involving novel view synthesis, this is not a critical consideration. However, by encouraging a unique Lambertian colour for each point in the scene where possible, we are able to better represent geometry by encouraging a view-independent appearance that must

be consistent with all other observations of the scene. In essence, our regularisation minimises the magnitude of the specular term with the assumption that the view-dependent appearance should be sparse through viewing angles.

## C. Total Variation of Unordered Samples on a Sphere

In this section, we provide some additional details regarding our total variation regularisation $\mathcal{R}_s$ on the specular colour $\mathbf{c}_s$. In Figure C.1 we provide a visual depiction of directions through which we sample candidate surface points in Surf-NeRF in 3D, and as coordinates on the surface of a sphere. To minimise total variation in specular colour, we select the $k$-nearest neighbours of a sample point on the surface of the hemisphere by angular distance. Selecting
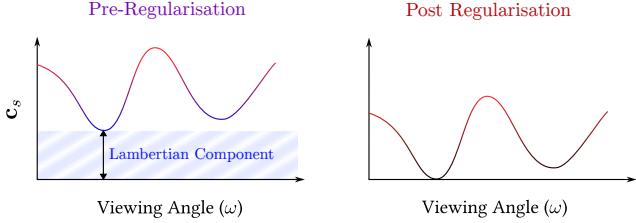
Figure B.1. Our motivation in regularising the view-dependent appearance for geometry within a NeRF. The specular term absorbs the Lambertian colour of geometry, increasing reliance on floaters and false geometry. Removing this Lambertian bias forces the model to place Lambertian geometry where it is able to, reducing reliance on floaters.
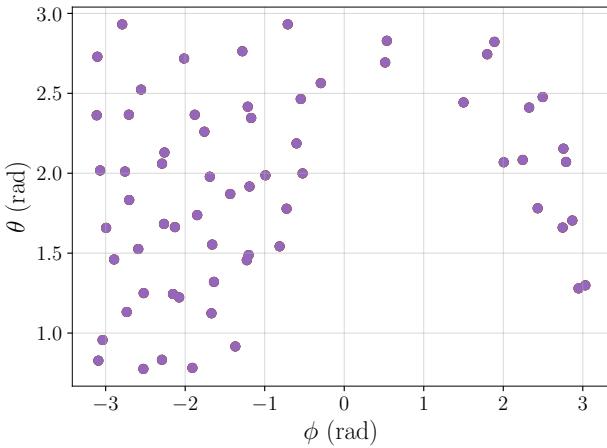


Figure C.1. Example directional polar coordinate map. The spliced forward and backward facing directions creates two interleaved sampling patterns, which is no longer regular. To address this, we use a graph total variation adapted for locations on the sphere accounting for irregularities.

nearest neighbours produces the unique property that the regularisation is continuous through all angles - that is, regularisation can occur at any point over the sampling sphere. Where samples are close to angle wrapping on the sphere, a regularisation loss is formed from the closest points irrespective of which hemisphere defined by $\phi$ they lie in.

We experimented with different numbers of $k$ when formulating this regularisation term. We found minimal improvement in geometric quality and marginally reduced visual fidelity going to higher numbers of $k$ above 3 most likely due to the small number of regularisation samples used and their comparative sparsity on the spherical domain.

## D. Model Details

In Figure D.1, we provide a visual depiction of our model structure incorporating the multiresolution hash encoding from ZipNeRF [6] and the physics-inspired structure of Ref-NeRF [36]. For our positionally encoded implementation, we use the same model parameters as those in the original Ref-NeRF paper.

From the base ZipNeRF implementation, we do not alter the hash encoding construction, keeping the two proposal and final NeRF networks at the same resolution. We retain the 512, 2048 and finally 8192 grid sizes in that order. We also maintain the proposed sampling strategy of the NeRF from ZipNeRF (64 samples for each proposal network, and 32 for the final NeRF), our experimentation did not indicate substantial changes to the final quality of the NeRF with the same 128 final samples used in Ref-NeRF. In particular, in our first surface assumption for Surf-NeRF we note that the majority of the radiance along a ray should be the result of only a handful of samples clustered at the surface; using a small number of final samples is in keeping with this assumption.

For the spatial MLP, we use a 2-layer MLP with 128 neurons in each layer, a deeper network compared to base ZipNeRF helping the NeRF to learn the additional channels under the Ref-NeRF parameterisation. We use a 3-layer, 256 neuron wide MLP to learn the specular colour similar to ZipNeRF. The respective components from Ref-NeRF, namely the integrated directional encoding (IDE), tone mapping reflection calculation and dot-product remain unchanged in their construction and addition to the network.

All weightings for losses in ZipNeRF [6] and Ref-NeRF [36] remain unchanged from previously published values. The Surf-NeRF regularisation terms have the following $\lambda$ weightings: our density smoothness $\mathcal{R}_d$ and normal consistency $\mathcal{R}_n$ terms have weightings $10^{-1}$, our specular bias term $\mathcal{R}_b$ has weighting $3 \cdot 10^{-2}$ and our specular total variation term $\mathcal{R}_s$ has weighting $10^{-3}$.

Our hash based implementations were trained for 25,000 iterations using a batch size of $2^{16}$, adjusted for GPU memory using the linear scaling rule. Our positionally encoded implementations were trained for 250,000 iterations using a batch size of $2^{14}$ rays. Optimisation followed the same learning schedule as ZipNeRF [6] and MipNeRF360 [4] for the hash-based and positionally encoded based implementations respectively.

## E. Surface Sampling Behaviour as Samples Converge

In this section we discuss briefly the behaviour of NeRF, and subsequently Surf-NeRF sampling as ray samples converge to a final value.

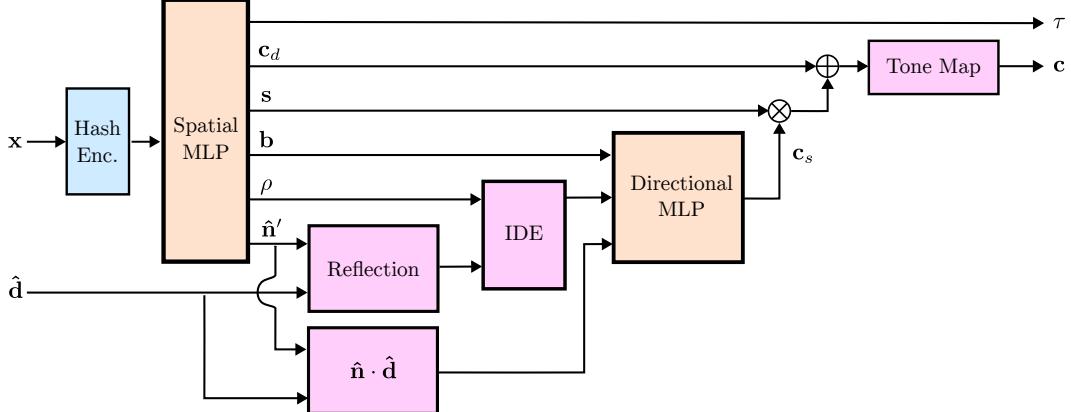Where the proposed sampling for the fine NeRF model

Figure D.1. Visual depiction of the adapted Ref-NeRF [36] model structure, leveraging the multi-resolution hash encoding in ZipNeRF [6]. We use the intermediate results for diffuse $\mathbf{c}_d$ and specular $\mathbf{c}_s$ colour terms to formulate the Surf-NeRF regularisation terms in terms of a Lambertian and specular scene appearance.
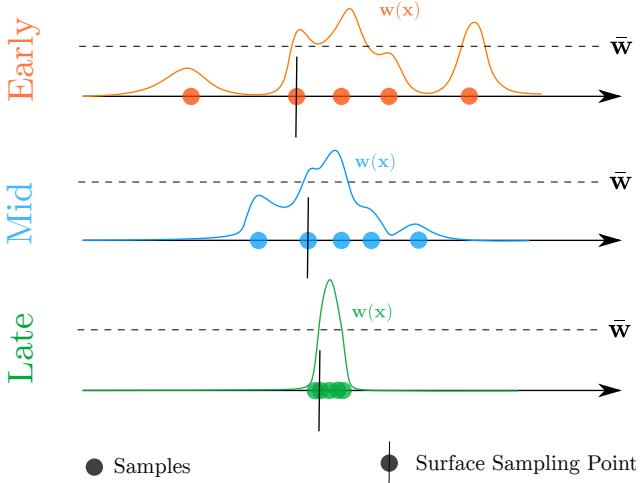


Figure E.1. The evolution of weights as the scene representation converges during training. Surf-NeRF uses a first surface assumption to determine the location of sampling, which is robust to outliers along the ray.

in MipNeRF360 produces highly clustered values, the spacing between subsequent samples becomes small; that is $\Delta t \to 0$. The ideal surface representation within a NeRF is in fact a delta function along the ray, where density is very high and clustered close to one another. An illustration of the evolution of the sample spacing is provided in Figure E.1.

In this limit, the Gaussians along a ray reach a minimum. Specifically, the sampling width $t_\delta$ approaches 0, whilst the mean sample approaches $t_i$,

$$t_\mu = \lim_{t_{i+1} \to t_i} \frac{t_{i+1} + t_i}{2} = t_i, \qquad (6a)$$

$$t_\delta = \lim_{t_{i+1} \to t_i} \frac{t_{i+1} - t_i}{2} = 0. \qquad (6b)$$

Evaluating the equations for mean and variance in MipN-eRF [3], we obtain,

$$\lim_{t_{i+1} \to t_i} \mu_t = t_\mu, \qquad (7a)$$

$$\lim_{t_{i+1} \to t_i} \sigma_t^2 = 0, \qquad (7b)$$

$$\lim_{t_{i+1} \to t_i} \sigma_r^2 = \frac{\dot{r}^2 t_\mu^2}{4}. \qquad (7c)$$

The samples turn from 3D anisotropic Gaussians, to planar isotropic Gaussians in the scene orientated along the ray, thereby sampling in the radial direction only.

In Surf-NeRF, we use the radial standard deviation $\sigma_r$ to scale our sampling sphere. As a result, our additional sampling scales to the certainty of a surface being present. When the samples are placed far apart as the NeRF is still converging to localised density along a ray, $t_\delta$ is nonzero and $\sigma_r$ is greater than $\dot{r}t/2$. Regularisation therefore occurs over a larger area, encouraging geometric smoothness and normal consistency to be based on samples further away, adding density, but ablating finer details. As training progresses, $\sigma_r$ approachs a minimum, ensuring sampling occurs within the bounds of a single pixel thereby retaining fine details and finetuning the representation.

## F. Koala Dataset Details

Here we provide some additional details regarding our captured Koala dataset. This dataset was captured using a Universal Robotics UR5e robotic arm, enabling ground truth poses to be captured and the same trajectory to be used between different captured scenes up to repeatability accuracy of the robot arm (less than one tenth of a mm). We

specifically curate scenes to have increased shape-radiance ambiguity by using objects with more complex reflection, compared to the existing *Shiny Real* [36] dataset, having large baselines between images thereby only sparsely sampling the directional radiance across view angles and capturing objects at several distances. Images were captured using a Basler acA1920-25uc camera, at constant gain and exposure in a controlled laboratory setting. This dataset provided a substantially more challenging benchmark given comparatively fewer images and more complex geometry.

# G. Per-Scene Results

In this section we present the per-scene results of our baseline and proposed methods across the *Shiny Objects*, *Shiny Real*, and *Koala* datasets. For the rendered *Shiny Objects* dataset we provide bold values demonstrating the improved geometric performance attained from Surf-NeRF in comparison to other baseline methods across both the traditional positional encoded variants and the hash based methods. For the real datasets we demonstrate comparable visual performance to the existing methods with minimal degradation compared to the Ref-NeRF baselines we benchmark against.

Table G.1 outlines results for the *Shiny Objects* [36] dataset. We show improved normals using our method across the positional encoded and grid-based approaches, and improved disparity RMSE scores. Figure G.1 provides additional qualitative results on this dataset showing separation of the Lambertian scene content and qualitatively improved geometry across all datasets. In Figure G.4 we also show the impact of our curriculum ablation on the coffee scene, a comparatively difficult scene owing to low texture and reflectively, demonstrating improved performance at extremely high levels of regularisation and the evolution of Lambertian and specular appearance with increasing frequency of regularisation.

Table G.2 demonstrates comparable visual fidelity of our proposed regularisation compared to the baseline approaches, with a small decrease overall owing to appearance and geometry being regularised in our curriculum learning framework. Figure G.2 demonstrates qualitative performance on this dataset, illustrating improved separation, depth and comparable view fidelity.

Finally, Table G.3 provides quantitative results for our captured dataset, showing aspects of improved visual fidelity in difficult scenes. These are illustrated accordingly in Figure G.3.

Table G.1. Per-Scene Results on the *Shiny Objects* [36] dataset. Yellow, orange, red indicate third, second and first performing scores. Bold indicates best score for model type in disparity RMSE. We achieve improvements in all cases for disparity, overall improved normals and comparable visual fidelity results over the dataset.

| | Model | PSNR | SSIM | MAE | RMSE |
|---|---|---|---|---|---|
| | MipNeRF | 26.99 | 0.921 | 47.44 | 0.127 |
| | MipNeRF+Diff | 26.55 | 0.920 | 53.49 | 0.130 |
| | Ref-NeRF | 31.17 | 0.957 | 15.49 | **0.123** |
| Car | Surf-NeRF (Ours) | 30.80 | 0.955 | 14.99 | **0.123** |
| | ZipNeRF | 27.20 | 0.930 | 23.49 | 0.207 |
| | Zip+Ref-NeRF | 30.37 | 0.955 | 16.37 | 0.205 |
| | Surf-NeRF (Ours) | 29.82 | 0.953 | 13.16 | **0.198** |
| | MipNeRF | 31.36 | 0.966 | 31.10 | 0.124 |
| | MipNeRF+Diff | 31.45 | 0.966 | 38.97 | 0.120 |
| | Ref-NeRF | 33.79 | 0.973 | 13.50 | **0.115** |
| Coffee | Surf-NeRF (Ours) | 33.56 | 0.972 | 12.54 | **0.115** |
| | ZipNeRF | 30.79 | 0.977 | 13.48 | 0.207 |
| | Zip+Ref-NeRF | 32.56 | 0.971 | 11.95 | 0.203 |
| | Surf-NeRF (Ours) | 31.99 | 0.968 | 14.79 | **0.185** |
| | MipNeRF | 29.03 | 0.943 | 75.28 | 0.117 |
| | MipNeRF+Diff | 27.95 | 0.934 | 72.18 | 0.112 |
| | Ref-NeRF | 29.31 | 0.952 | 40.96 | **0.109** |
| Helmet | Surf-NeRF (Ours) | 30.09 | 0.958 | 26.68 | **0.109** |
| | ZipNeRF | 27.07 | 0.945 | 24.96 | 0.223 |
| | Zip+Ref-NeRF | 32.37 | 0.975 | 10.42 | 0.193 |
| | Surf-NeRF (Ours) | 32.93 | 0.974 | 8.91 | **0.192** |
| | MipNeRF | 45.90 | 0.996 | 67.07 | 0.152 |
| | MipNeRF+Diff | 45.62 | 0.996 | 69.10 | 0.145 |
| | Ref-NeRF | 46.24 | 0.997 | 15.82 | **0.136** |
| Teapot | Surf-NeRF (Ours) | 45.76 | 0.996 | 16.07 | **0.136** |
| | ZipNeRF | 46.11 | 0.997 | 9.10 | 0.192 |
| | Zip+Ref-NeRF | 46.09 | 0.997 | 7.62 | **0.157** |
| | Surf-NeRF (Ours) | 44.23 | 0.996 | 6.17 | **0.157** |
| | MipNeRF | 23.17 | 0.890 | 59.97 | 0.104 |
| | MipNeRF+Diff | 22.63 | 0.864 | 66.27 | 0.104 |
| | Ref-NeRF | 25.52 | 0.977 | 43.68 | **0.102** |
| Toaster | Surf-NeRF (Ours) | 24.76 | 0.911 | 41.20 | **0.102** |
| | ZipNeRF | 23.91 | 0.911 | 26.33 | 0.216 |
| | Zip+Ref-NeRF | 24.29 | 0.921 | 27.16 | 0.215 |
| | Surf-NeRF (Ours) | 21.67 | 0.879 | 23.72 | **0.198** |

Table G.2. *Shiny Real* [36] dataset results. Yellow, orange, red indicate third, second and first performing scores. Applying our regularisation does not significantly degrade visual fidelity compared to unregularised variants.

| Scene | Model | PSNR | SSIM |
|---|---|---|---|
| | MipNeRF | 22.37 | 0.452 |
| | MipNeRF+Diff | 22.34 | 0.460 |
| | Ref-NeRF | 22.17 | 0.413 |
| Garden Spheres | Surf-NeRF (Ours) | 22.33 | 0.443 |
| | ZipNeRF | 21.68 | 0.539 |
| | Zip+Ref-NeRF | 21.69 | 0.540 |
| | Surf-NeRF (Ours) | 20.51 | 0.486 |
| | MipNeRF | 24.71 | 0.600 |
| | MipNeRF+Diff | 24.51 | 0.594 |
| | Ref-NeRF | 24.80 | 0.602 |
| Sedan | Surf-NeRF (Ours) | 24.72 | 0.614 |
| | ZipNeRF | 25.99 | 0.730 |
| | Zip+Ref-NeRF | 23.92 | 0.687 |
| | Surf-NeRF (Ours) | 23.71 | 0.684 |
| | MipNeRF | 24.09 | 0.576 |
| | MipNeRF+Diff | 24.00 | 0.576 |
| | Ref-NeRF | 23.99 | 0.572 |
| Toycar | Surf-NeRF (Ours) | 24.06 | 0.571 |
| | ZipNeRF | 23.23 | 0.601 |
| | Zip+Ref-NeRF | 23.24 | 0.609 |
| | Surf-NeRF (Ours) | 23.20 | 0.620 |

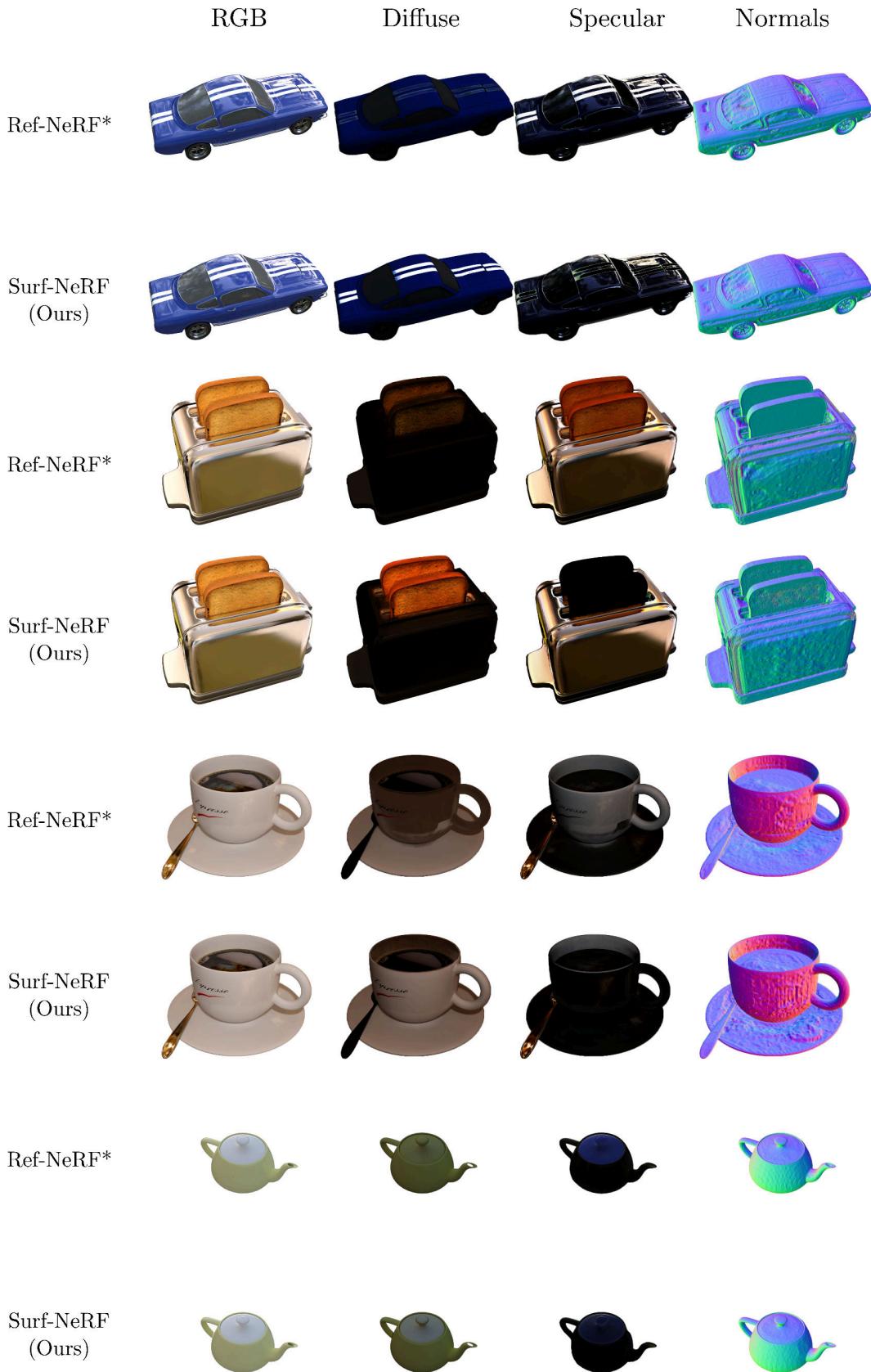|  | RGB | Diffuse | Specular | Normals |
|---|---|---|---|---|

Figure G.1. Examples from the *Shiny Objects* dataset [36] for the hash based variants. Diffuse components are correctly removed from the specular term of the rendering, and normals show more complete surface structure. Ref-NeRF* indicates the Zip+Ref-NeRF baseline.
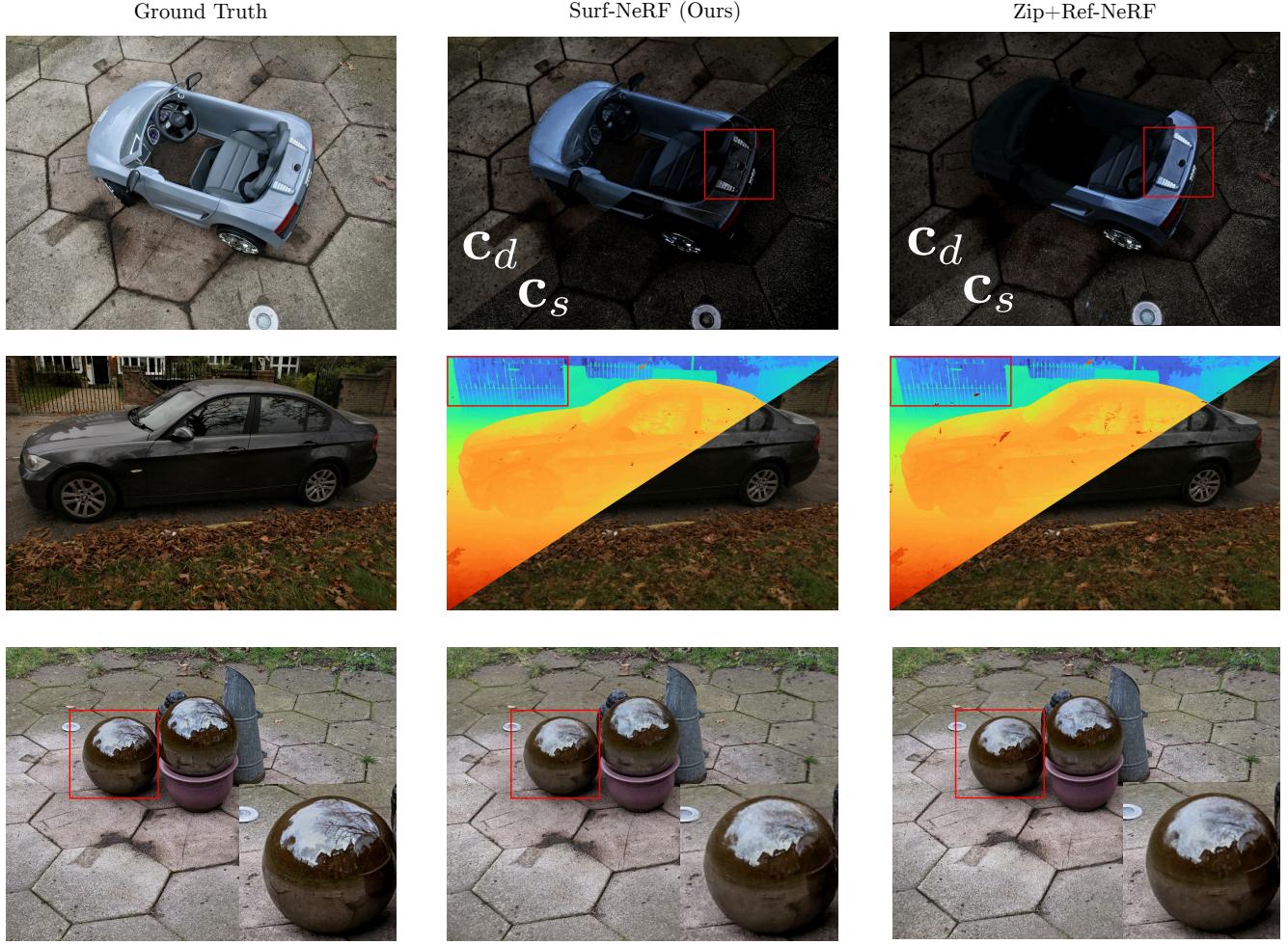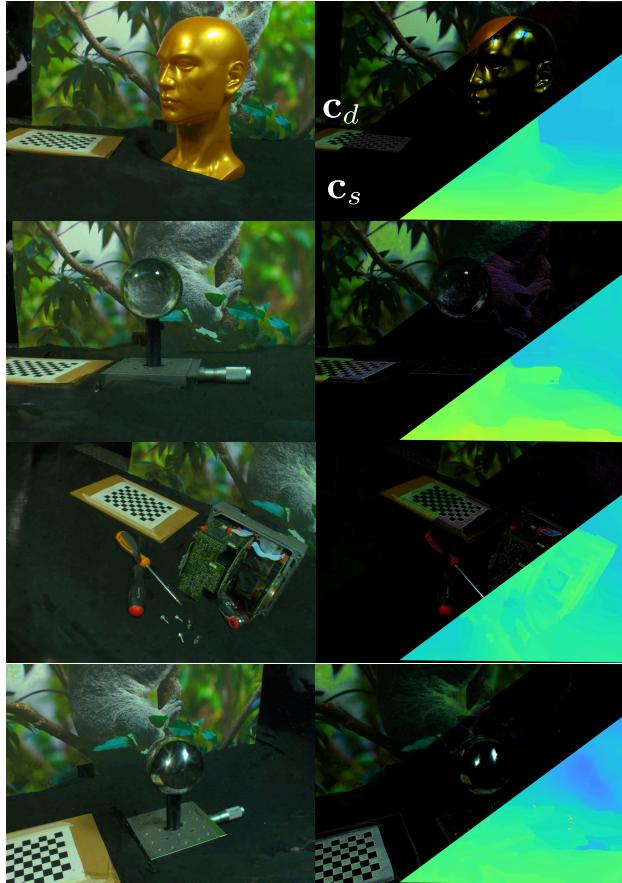
Figure G.2. Results on the *Shiny Real* dataset [36]. Surf-NeRF attains separation of diffuse and specular components $\mathbf{c}_d$ and $\mathbf{c}_s$ (top), a decreased reliance on floaters, more consistent depths, and qualitatively better representation of background content (middle), and comparable visual fidelity of reflective objects to other reflection parameterised NeRF variants.

Table G.3. Surf-NeRF results on our captured *Koala* dataset. Bold indicates best performing scores. Regularisation helps the reflection-parameterised model around complex geometry and appearance on this dataset.

| Scene | Model | PSNR | SSIM |
|---|---|---|---|
| Gold Head | ZipNeRF | **29.00** | 0.680 |
| | Zip+Ref-NeRF | 28.34 | **0.681** |
| | Surf-NeRF (Ours) | 27.62 | 0.680 |
| Shiny Ball | ZipNeRF | **28.59** | 0.641 |
| | Zip+Ref-NeRF | 15.47 | 0.435 |
| | Surf-NeRF (Ours) | 28.17 | **0.646** |
| Lidar Guts | ZipNeRF | 27.42 | 0.664 |
| | Zip+Ref-NeRF | **27.55** | 0.665 |
| | Surf-NeRF (Ours) | 26.84 | **0.666** |
| Crystal Capture | ZipNeRF | 26.14 | **0.665** |
| | Zip+Ref-NeRF | 23.40 | 0.577 |
| | Surf-NeRF (Ours) | **26.33** | 0.627 |

Surf-NeRF (Ours)                          Zip+Ref-NeRF



$\mathbf{c}_d$
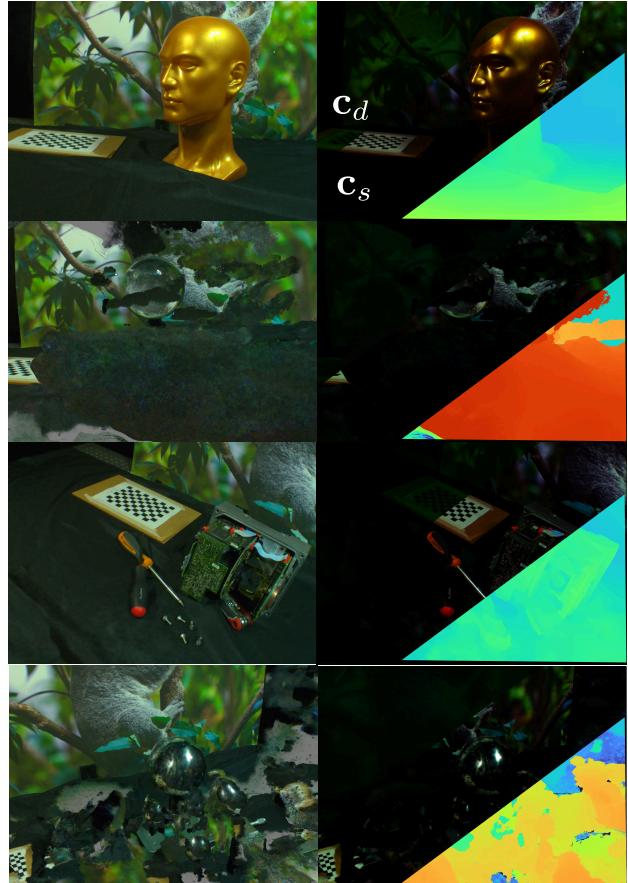
$\mathbf{c}_s$

$\mathbf{c}_d$

$\mathbf{c}_s$

Figure G.3. Our captured *Koala* dataset. We show improved separation of scene appearance between view-dependent and -independent content, more stable performance for scene content with high curvature and high view dependence. Renderings from models are on the left frame, with diffuse, specular and depth renderings on the right.

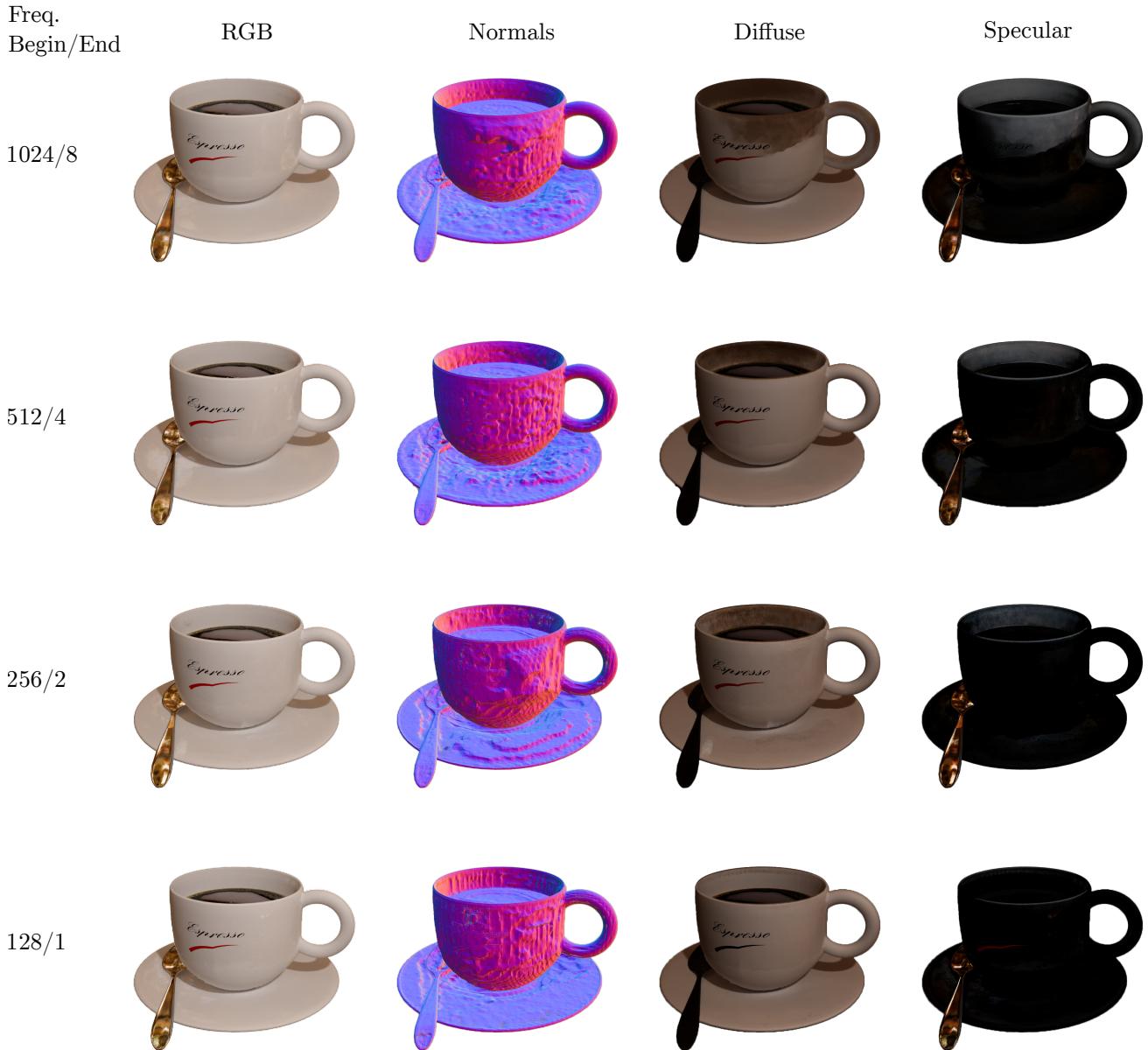| Freq. Begin/End | RGB | Normals | Diffuse | Specular |
|---|---|---|---|---|
| 1024/8 | | | | |
| 512/4 | | | | |
| 256/2 | | | | |
| 128/1 | | | | |

Figure G.4. Visualisation of the curriculum learning ablation study in the main text. As frequency increases, we see geometry become more consistent and appearance separate more readily. We found that there was tradeoffs to how quickly we could regularise - as density is changed, the model must alter appearance to "take-up" the changes to geometry. The 512/4 frequency provided a medium ground between strength of geometric improvement and the model's ability to alter appearance for the hash-based methods.