

TraM-NeRF: Tracing Mirror and Near-Perfect Specular Reflections through Neural Radiance Fields

Leif Van Holland, Ruben Bliersbach, Jan U. Müller, Patrick Stotko, Reinhard Klein

University of Bonn



Figure 1: Examples of novel views rendered with our proposed approach on scenes with mirror surfaces (left, center) and near-perfect specular surfaces (right).

Abstract

Implicit representations like Neural Radiance Fields (NeRF) showed impressive results for photorealistic rendering of complex scenes with fine details. However, ideal or near-perfectly specular reflecting objects such as mirrors, which are often encountered in various indoor scenes, impose ambiguities and inconsistencies in the representation of the reconstructed scene leading to severe artifacts in the synthesized renderings. In this paper, we present a novel reflection tracing method tailored for the involved volume rendering within NeRF that takes these mirror-like objects into account while avoiding the cost of straightforward but expensive extensions through standard path tracing. By explicitly modeling the reflection behavior using physically plausible materials and estimating the reflected radiance with Monte-Carlo methods within the volume rendering formulation, we derive efficient strategies for importance sampling and the transmittance computation along rays from only few samples. We show that our novel method enables the training of consistent representations of such challenging scenes and achieves superior results in comparison to previous state-of-the-art approaches.

CCS Concepts

• **Computing methodologies** → Image-based rendering; Ray tracing; Reflectance modeling;

1. Introduction

3D reconstruction and modeling of real-world scenes has been a major research field for decades and plays a crucial role in a diverse range of applications such as video gaming, movies, advertisement, education as well as AR and VR scenarios. With the recent emergence of neural scene representations and, especially, Neural Radiance Fields (NeRF) [MST*20], a compelling degree of photorealism and immersion of the rendered views has been

achieved which inspired many further developments [ZRSK20; RPLG21; BMT*21; MESK22; WWG*22; CZL*22]. By combining graphics-based volume rendering with an efficient representation of scene density and radiance using neural networks in terms of multilayer perceptrons (MLP), NeRF enables capturing various effects including view-dependent changes of object appearances or volumetric phenomena like clouds. However, objects with ideal and near-perfect specular reflection behavior which are often encountered in various scenarios and, in particular, many indoor scenes

impose a significant challenge to the representation capabilities of radiance fields as they induce a very specific pattern in the light transport. For the case of a planar mirror, a symmetric version of the visible scene parts can be observed which appears to be located behind the mirror and gives the illusion of viewing the respective content through a window. While, a-priori, this ambiguity results into two plausible and consistent interpretations of the structure of the surrounding environment, additional views directly from behind the mirror object allows to resolve the scenario as the representation of the virtual mirrored scene collides with the observations. As a consequence, severe artifacts will be introduced in the scene representation of NeRF as the underlying volume rendering approach for rendering traces the primary viewing rays and, in turn, implicitly always prefers the inconsistent interpretation. Several approaches addressed this issue by decomposing the scene into two or more individually consistent radiance fields [GKB*22; YQCR23] or employ standard path tracing [ZXY*23; MVKF23] in combination with an extended volumetric field to infer normal directions and specular reflection probabilities [ZBC*23]. However, this significantly increases the computational burden both in terms of training performance as well as rendering speed and limits their application into other sophisticated and advanced NeRF approaches.

In this paper, we direct our attention towards an efficient formulation of reflection tracing within the volume rendering procedure of NeRF that can be easily adopted in several NeRF variants to enhance their capabilities in handling mirror-like objects. To this end, we extend the single-ray absorption volume integration by considering the contributions of reflected radiance towards the observed radiance by the camera, effectively moving our model closer to full physically interpretable light transport in the process. Our proposed method, referred to as TraM-NeRF, explicitly integrates reflected radiance by first annotating near-specular surfaces. Combining NeRF volume rendering and ray-tracing with physically plausible materials at intersection points introduces an inductive bias into the training of TraM-NeRF that enables it to learn a single coherent scene representation, even when geometry has only been observed in a reflection. Our combined radiance estimator allows us to reduce its variance compared to a standard Monte-Carlo approach without additional computational overhead by decreasing the time spend on transmittance computation along rays.

In summary, the key contributions of this work are:

- We present TraM-NeRF, an architecture-agnostic extension of NeRF that efficiently represents scenes with mirror-like surfaces, modeling high-frequency reflections in a physically plausible manner within a single coherent scene representation.
- We derive a transmittance-aware formulation of the rendering equation to explicitly model reflected radiance at mirror-like surfaces. Additionally, we introduce efficient strategies for importance sampling and transmittance computation, resulting in a reduction in the number of network evaluations compared to Monte-Carlo estimation.
- We demonstrate the benefits of our formulation in comparison to previous state-of-the-art methods on a variety of challenging scenes some of which include multiple mirror-like surfaces.

The code of our implementation is available at <https://github.com/Rubikalubi/TraM-NeRF>.

2. Related Work

2.1. Neural Scene Representations

Synthesizing novel views of complex scenes has gained increasing interest due to the promising results achieved with neural scene representations [LSS*19; SZW19; NMOG20; BXS*20a; BXS*20b]. Among these, especially the work on Neural Radiance Fields [MST*20] excels in terms of the quality and degree of photorealism of the rendered images and has become very popular, also due to its simple but effective formulation. In particular, NeRF leverages volume rendering to accumulate the scattered lighting contributions along the traced viewing rays which are represented using volumetric density and view-dependent radiance and parametrized using MLPs. Various extensions have been developed to further enhance the performance and quality of the original approach such as accelerating the training [MESK22; CXG*22; FYT*22] as well as the rendering processes [RPLG21; GJK*21], reducing aliasing artifacts by replacing ray-based marching with an integration of 3D conical frustums [BMT*21; BMV*22], rendering fine details at very high-resolution [WWG*22; WLS*22; LLGG23], or lifting its capabilities to also handle unbounded scenes [ZRSK20; BMV*22] and to reconstruct from in-the-wild image collections [MRS*21; CZL*22; FMW*23] or low dynamic range images with low or varying exposure [HZF*22; MHM*22].

Besides these advances, the underlying representation given by volumetrically baked radiance and density does not account for manipulation tasks like exchanging the environment illumination, so significant effort has been spent into more plausible, physics-inspired scene representations. Thus, various methods [ZSD*21; BBJ*21; SDZ*21; BEK*22; JLY*23] considered factorizing the radiance field into shape with normals, surface material parameters in terms of a Bidirectional Reflectance Distribution Function (BRDF) as well as environment illumination. Further approaches [ZLW*21; FSV*23; LCL*23; WHL*23; GHZ*23] replaced the density-based shape representation by implicit surfaces via signed distance functions (SDF) for a more accurate estimation of the object geometry and normals.

2.2. Specular Reflections in Neural Representations

Objects with highly reflective materials often exhibited in captured scenes impose are challenging to reconstruct in decomposed representations of neural radiance fields and have, in turn, attracted increasing attention. Ref-NeRF [VHM*22] reparametrizes the observed radiance based on the local normal vector and its angle to the view direction to a simpler model that shares common structures across multiple views. PhySG [ZLW*21] employs Spherical Gaussians to represent specular reflections in the BRDF which has been later extended by splitting the illumination into a direct and indirect component each modeling an individual specular reflection [ZSH*22]. Ref-NeuS [GHZ*23] detect anomalies in the rendered images caused by reflections and incorporate a respective reflection score into the photometric loss as a guidance. Other works instead directly trace reflections either only in the ideal reflection direction assuming a low material roughness [LCL*23] or using path tracing evaluated via Monte-Carlo estimators [WHL*23]. Recently, volumetric microflake [ZXY*23]

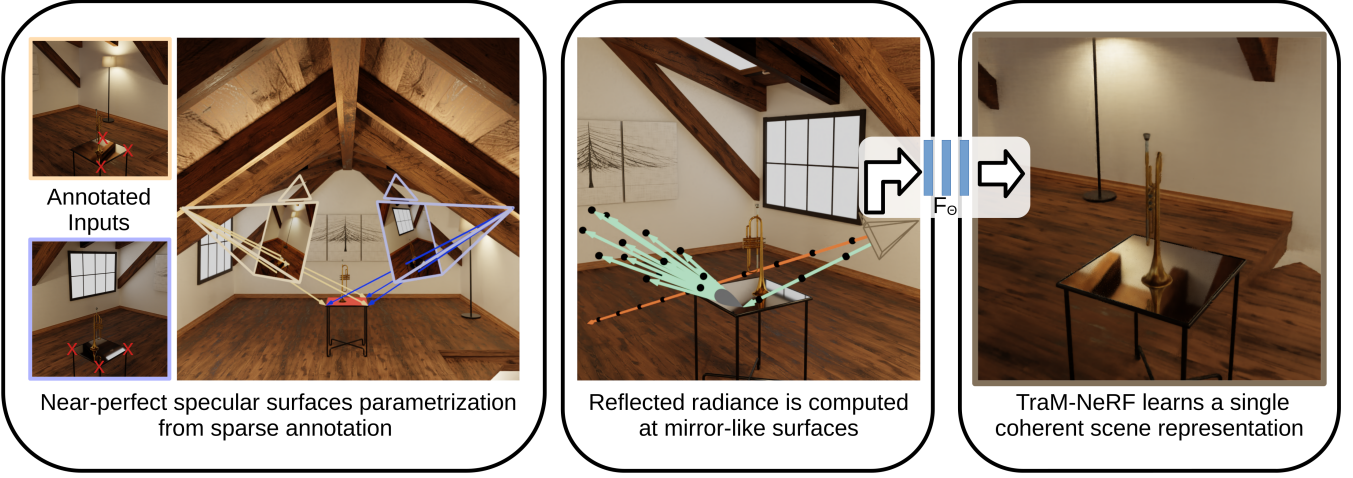


Figure 2: Overview of our proposed method TraM-NeRF. Our approach to parameterize nearly-specular surfaces using only sparse annotations. We introduce a radiance estimator, a crucial component of TraM-NeRF, which combines volume and reflected radiance integration for training and rendering the model. TraM-NeRF learns to represent the observed radiance in a single coherent network.

and microfacet [MVKF23] fields presented a hybrid rendering approach by combining the ray marching of volume rendering with importance-sampled path tracing according to the distribution of the micro structures.

Most closely related to our work are techniques that explicitly model mirror reflections within the scene. NeRFReN [GKB*22] decomposes the scene into two independently traversed and rendered radiance fields consisting of an ordinary NeRF for the transmitted radiance as well as an additional NeRF that only covers the reflected radiance. The final synthesized image is then obtained by blending the results for the transmitted and reflected fields. MS-NeRF [YQCR23] learns radiance and weights into multiple feature fields that are decoded by small MLPs for rendering and then blended together. Mirror-NeRF [ZBC*23] follows a different direction by representing the scene in a single radiance field and instead further tracing the rays in the ideal reflection direction after hitting a mirror. The respective normal directions and reflection probabilities used for reflecting the rays are additionally learned in the volumetric neural field similar as in the microflake/microfacet fields [ZXY*23; MVKF23].

In contrast to the aforementioned approaches, the primary focus of this work lies in the efficient rendering of unified radiance fields of scenes with mirror and near-perfect specular reflecting objects using only a low number of network evaluations for each importance-sampled reflection ray while achieving a significantly lower variance than standard Monte-Carlo estimators.

3. Method

To this end, we start with a brief overview of NeRF. Next, we introduce our radiance estimator, a crucial component of TraM-NeRF, which combines volume and reflected radiance integration for rendering and model training. We then discuss our approach to parameterize nearly-specular surfaces using sparse annotations. Finally,

we provide implementation and training details for transparency and reproducibility.

3.1. Neural Radiance Fields

We build upon the neural implicit scene representation proposed by MILDENHALL et al. [MST*20] which uses a simple MLP F_Θ to infer a RGB color value $c \in \mathbb{R}^3$ and a density $\sigma \in \mathbb{R}$ for a given spatial location $x \in \mathbb{R}^3$ and viewing direction $d \in \mathbb{R}^3$. In order to also capture high-frequency details, x is first lifted into a higher-dimensional space using a positional encoding

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

To render an image using F_Θ , we define a camera and consider rays starting at the camera center $o \in \mathbb{R}^3$ with direction $d \in \mathbb{R}^3$, such that $r(t) = o + t d$ represents a point on the ray. In the original NeRF formulation, the observed radiance

$$C_e(r) = \int_{t_n}^{t_f} T(t_n \rightarrow t) \sigma(r(t)) c(r(t), d) dt \quad (2)$$

corresponding to ray r is given by the volume integration through an absorbing medium where $T(t_n \rightarrow t) = \exp(-\int_{t_n}^t \sigma(r(s)) ds)$ represents the accumulated transmittance up to distance t , and $t_n, t_f \in \mathbb{R}$ describe the near and far plane of r . To be able to compute eq. (2) in practice, the integral is numerically approximated using quadrature [Max95] at locations t_k along the ray, yielding the following discrete sum

$$C_e(r) \approx \hat{C}_e(r) = \sum_{k=1}^K T_k \left(1 - e^{-\sigma_k \delta_k} \right) c_k. \quad (3)$$

Here, $\delta_k = t_k - t_{k-1}$ is the distance between successive locations and $T(t_n \rightarrow t_k) \approx T_k = e^{-\sum_{j=1}^{k-1} \sigma_j \delta_j}$ approximates the accumulated transmittance. To prevent using only a discrete subset of locations,

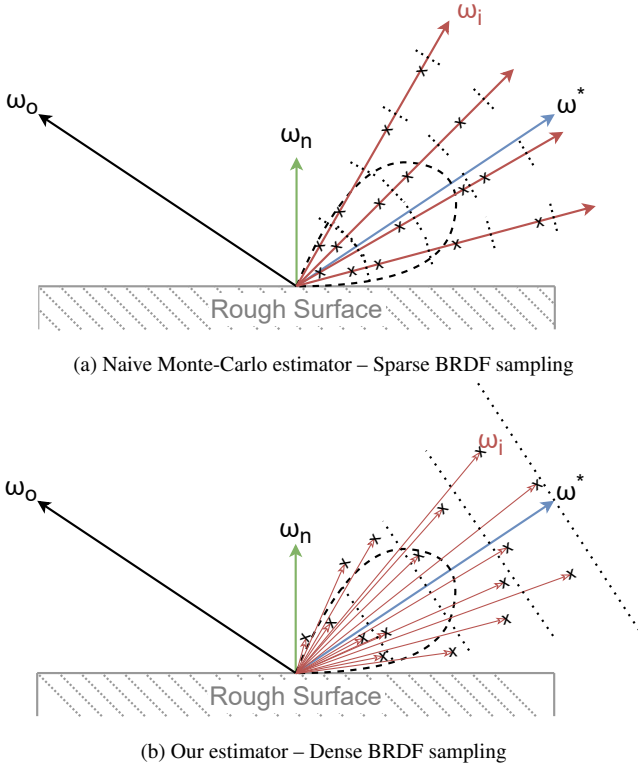


Figure 3: Patterns for BRDF sampling and network evaluation, resulting in different radiance estimators. In the standard Monte Carlo approach, the network is evaluate at positions (indicated by cross markers) chosen using stratified sampling for each sampled direction ω_i . (b) Our estimator draws directional samples within segments (dashed lines) along the ideal reflection direction ω^* resulting in a higher angular coverage of the specular lobe with the same number of network evaluations.

we use stratified sampling for the t_k , as proposed by MILDENHALL et al. [MST*20].

During training, the parameters Θ of the MLP are optimized via gradient descent using a photometric loss \mathcal{L} defined as the mean squared error between the ground-truth colors $C^*(r)$ and the rendered images $\hat{C}_e(r)$ over a batch of rays R :

$$\mathcal{L} = \frac{1}{|R|} \sum_{r \in R} \|C^*(r) - \hat{C}_e(r)\|_2^2. \quad (4)$$

3.2. Radiance Integration at Near-Perfect Specular Surfaces

Assume that a camera ray $r(t) = o_c + t d_c$ intersects with the near-specular surfaces which have been detected in a point x . To allow a model to learn a consistent representation of observed radiance in a single radiance field, we drop the assumption of NeRF that the ray terminates (i.e. the transmittance vanishes) at an opaque surface. Instead, TraM-NeRF relies on the rendering equation [Kaj86] to compute its predicted radiance at intersection points, which states that the radiance $C(x, \omega_o)$ at a point x when observed from direction

ω_o is the sum of the emitted radiance $C_e(x, \omega_o)$ and the reflected radiance $C_r(x, \omega_o)$:

$$C(x, \omega_o) = C_e(x, \omega_o) + C_r(x, \omega_o). \quad (5)$$

Here, the reflected radiance is obtained by evaluating the transport integral

$$C_r(x, \omega_o) = \int_{\Omega} f(x, \omega_i, \omega_o) C(x, \omega_i) \cos \theta_i d\omega_i \quad (6)$$

over the visible hemisphere Ω where $f(x, \omega_i, \omega_o)$ denotes the BRDF and θ_i is the angle between the surface normal at x and the direction of incoming light ω_i . Note that by convention the direction of outgoing light ω_o faces outwards from x . Consequently, we set the direction of outgoing light to be $\omega_o = -d_c$.

Combining Surface Rendering and Volume Integration In order to combine the radiance integration and volume integration, TraM-NeRF assumes that a primary ray scatters into multiple reflected rays at the intersecting point. Since this ray has passed through an absorbing medium, the combined radiance of the out-branching rays should be attenuated by the transmittance along the intersecting ray. Whereas NeRF integrates the density along a ray from a starting point close to the camera position up to a point which is chosen a-priori based on the extent of the scene, TraM-NeRF modifies the upper integration bound to stop at the intersection point. We formalize this concept by introducing a ray length function $\tau(x, \omega)$ which returns the length from the ray origin to the point where it intersects with the detected geometry. Therefore, we obtain a transmittance-aware version of the rendering equation

$$C(x, \omega_o) = C_e(x, \omega_o) + T_{\omega_o}(t_n \rightarrow \tau(x, \omega_o)) \cdot C_r(x, \omega_o). \quad (7)$$

which takes the attenuation from the absorbing medium into account by multiplying the reflected radiance with the transmittance $T_{\omega_o}(t_n \rightarrow \tau(x, \omega_o))$. The emitted radiance observed at point x from direction ω_o is computed following [MST*20] by raymarching through the emissive volume until the intersection point is reached:

$$C_e(x, \omega_o) = \int_{t_n}^{\tau(x, \omega_o)} T_{\omega_o}(t_n \rightarrow t) \sigma_{\omega_o}(t) c_{\omega_o}(t) dt. \quad (8)$$

Note, our modified version retains the offset t_n to prevent double-counting the intersection point in emitted and reflected radiance calculations.

Monte-Carlo Estimator of Reflected Radiance Our efficient reflected radiance estimator builds upon an established approximation method for the transport integral, employing importance sampling to evaluate a Monte-Carlo estimator. This estimator is then modified to reduce additional variance introduced when importance sampling a BRDF function, all while keeping the number of network evaluations constant. These adjustments enhance the computational efficiency in determining the reflected radiance.

Considering the transport integral in eq. (6), the respective estimator sampling from a candidate distribution $p(\cdot)$ is given by

$$C_r(x, \omega_o) \approx \frac{1}{n} \sum_{i=1}^n \frac{f(x, \omega_o, \omega_i)}{p(\omega_i)} C(x, \omega_i) \quad (9)$$

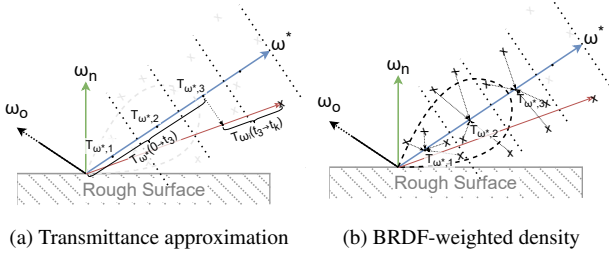


Figure 4: Transmittance approximation used by our estimator. (a) The transmittance $T_{\omega_i}(t_n \rightarrow t_k)$ in direction ω_i is approximated using the transmittance $T_{\omega^*}(t_n \rightarrow t_{k-1})$ along the ideal reflection direction ω^* for near-perfect specular surfaces. (b) The transmittance in the ideal reflection direction is computed using the BRDF-weighted density of samples within each segment.

where we make use of a convenient overlap in notation and re-define ω_i to be the i -th directional sample drawn from Ω with probability $p(\omega_i)$.

In order to obtain a suitable candidate distribution derived from the BRDF f , we utilize the well-established microfacet theory to model f at the intersection point, assuming roughness arises from a height field of tiny facets distributed according to a distribution $D_\alpha(\cdot)$ with roughness parameter α [CT82]. In particular, we use the widely used GGX reflection model [WMLT07] which defines the BRDF to be

$$f(x, \omega_o, \omega_i) = \frac{F(\omega_i, h) G(\omega_o, \omega_i, h) D_\alpha(h)}{4 \cos \theta_i \cos \theta_o} \quad (10)$$

where h is the half-vector between the incoming and outgoing direction, θ_o the angle between ω_o and the surface normal, F is the Fresnel term, and $G(\cdot)$ a coefficient describing the average attenuation that results from shadowing and masking between microfacets. For formal definitions of D_α , F , and G , please refer to [WMLT07]. Utilizing an optics-based analytical GGX reflectance model ensures physical consistency in the learned radiance function and comes with the additional benefit of a well-studied importance sampling technique [Hei18; DB23]. We use visible normal sampling (VNDF) which defines a candidate distribution

$$p(h) = \frac{\max\{0, \langle \omega_o | h \rangle\} G_1(\omega_o) D_\alpha(h)}{4 \cos \theta_i \cos \theta_o} \quad (11)$$

that takes the average attenuation due to microfacet masking G_1 into consideration and has a closed-form sampling routine [Hei18]. Note that VNDF is a distribution over the half-vectors h instead of incoming light directions. However, the incoming light direction can be computed via a reflection of the outgoing direction about the half-vector. Thus, the resulting estimator for the reflected radiance is

$$C_r(x, \omega_o) \approx \frac{1}{n} \sum_{i=1}^n \underbrace{\frac{F(\omega_i, h) G(\omega_o, \omega_i, h)}{G_1(\omega_o)}}_{=: f'(\omega_o, \omega_i)} C(x, \omega_i). \quad (12)$$

Efficient Reflected Radiance Approximation In order to motivate our efficient radiance approximation, we first review the num-

ber of network evaluations required by the estimator when the volume integral is discretized using the same assumption made in NeRF [MST*20], that is assuming that the radiance along the ray direction is piece-wise constant. To simplify both the argument and notation and without loss of generality, we assume that only the camera ray intersects with any detected planar surface. This simplification combined with the discretized volume integration yields the following reflected radiance estimator:

$$C_r(x, \omega_o) = \frac{1}{n} \sum_{i=1}^n f'(\omega_o, \omega_i) \sum_{k=1}^K T_{\omega_i}(t_n \rightarrow t_k) \left(1 - e^{-\sigma_{\omega_i, k} \delta_{\omega_i, k}}\right) c_{\omega_i, k} \quad (13)$$

Figure 3a provides a visualization of the network evaluation pattern for each directional sample. A crucial observation is that computing this equation entails K network evaluations for each direction sampled using VNDF. To achieve a radiance estimate with minimal noise, a large number of directional samples are required, necessitating a significant number of costly network evaluations.

In light of the aforementioned challenges, we aim to improve the computational efficiency of this procedure in TraM-NeRF. Our strategy involves increasing the number of directional samples without the need for additional network evaluations. This optimization leverages the observation that scenes with diffuse or low frequency surfaces reflections can be adequately handled by the standard NeRF model. However, it encounters difficulties in representing scenes with surfaces which display high-frequency reflections that vary significantly with the viewing direction. In TraM-NeRF, we combine this observation with the assumptions of scene bound-ness and locally smooth network-predicted density. These assumptions allows our estimator to focus primarily on estimating reflected radiance for near-specular surfaces. Consequently, our estimator can assume a narrow spread of light directions in the samples. By combining these insights, we expect that the transmittance remains nearly constant with respect to the sampled directions. In particular, we approximate the transmittance along all sampled directions ω_i with the transmittance along the ideal reflection direction ω^* :

$$T_{\omega_i}(t_n \rightarrow t_k) \approx T_{\omega^*}(t_n \rightarrow t_k). \quad (14)$$

This way, we can interchange the order of the Monte Carlo integration and the NeRF volume integration:

$$\begin{aligned} C_r(x, \omega_o) &= \frac{1}{n} \sum_{i=1}^n f'(\omega_o, \omega_i) \sum_{k=1}^K T_{\omega_i}(t_n \rightarrow t_k) \left(1 - e^{-\sigma_{\omega_i, k} \delta_{\omega_i, k}}\right) c_{\omega_i, k} \\ &\approx \sum_{k=1}^K T_{\omega^*}(t_n \rightarrow t_k) \frac{1}{n} \sum_{i=1}^n f'(\omega_o, \omega_i) \left(1 - e^{-\sigma_{\omega_i, k} \delta_{\omega_i, k}}\right) c_{\omega_i, k} \end{aligned} \quad (15)$$

Intuitively, our estimator divides the ideal reflection ray into K segments and traces transmittance solely along the ideal reflection direction. Within each segment, we randomly select positions and evaluate n directional samples, each contributing to the calculation only once with their impact attenuated based on the transmittance along the ideal reflection direction, which is depicted in Figure 4a. These positions for evaluating the directional samples are uniformly chosen from the interval $[\frac{t_{k-1}+t_k}{2}, \frac{t_k+t_{k+1}}{2}]$. A visualization of this sampling strategy and its resulting evaluation points is

shown in Figure 3b. For a more in-depth explanation and its adaptation to a hierarchical optimization procedure, please refer to Section 3.4.

Our estimator computes transmittance once per segment, enabling us to increase the number of directional samples without requiring additional network evaluations to accumulate transmittance. This trade-off balances transmittance precision against directional sample count, reducing noise in reflected radiance. To further reduce the number of network evaluations, we use an average of BRDF-weighted density predictions per segment to calculate the transmittance

$$T_{\omega^*}(t_n \rightarrow t_k) \approx e^{-\sum_{j=1}^{k-1} \left[\frac{1}{n} \sum_{i=1}^n f'(\omega_o, \omega_i) \sigma_{\omega_i, j} \right] \delta_j}, \quad (16)$$

which is visualized in Figure 4b.

3.3. Mirror Parameterization and Annotation

TraM-NeRF infers the position of near-specular planar surfaces from sparsely annotated data, relying only on annotations of corner points marked in a few images. In practice, we get sufficiently accurate annotations using only two annotated input images per scene.

We represent planar surfaces as triplets of triangle vertices $T = (v_1, v_2, v_3)$, $v_i \in \mathbb{R}^3$ which allows for efficient intersection tests with rays in the rendering step [MT97]. Given the screen space annotations of three corners in at least two images and their camera poses, the annotations correspond to rays through the scene. In particular, the j -th annotation of vertex v_i defines a ray $r_{ij}(t) = o_j + t d_{ij}$ with camera origin o_j and ray direction d_{ij} with $\|d_{ij}\|_2 = 1$. The estimated 3D location \hat{v}_i of the vertex is then given as the point minimizing the lengths of the orthogonal projection onto each ray:

$$\hat{v}_i = \min_v \sum_j \|v - (o_j + d_{ij}(v - o_j) d_{ij})\|_2^2. \quad (17)$$

By limiting TraM-NeRF to planar surfaces, we can additionally exploit the property that all vertices of a mirror lie on the same plane. To increase the robustness against inaccuracies in the annotations, we compute the normal of that plane using principal component analysis applied on the set of annotated vertices of the planar surface.

3.4. Implementation and Training Details

To assess and compare our estimator for reflected radiance, TraM-NeRF leverages the NeRF framework [MST*20] and is implemented using PyTorch [PGM*19]. We chose to use the standard NeRF implementation to ensure a clear comparison of the improvements resulting from our contributions and to avoid potential confusion in the assessment of enhancements attributable specifically to our estimator in comparison to those resulting from different unrelated improvements. Nevertheless, our estimator is adaptable to various implementations of the radiance field networks, making it compatible with methods that enhance parametrization [BMT*21; BMV*22] or architecture [MESK22; CXG*22]. For training TraM-NeRF, we use a modified version of the NeRF training protocol [MST*20], using the Adam optimizer with a learning rate of 10^{-3}

without decay. The Adam hyperparameters remained at their default values: $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\epsilon = 10^{-7}$. The model underwent 6×10^4 iterations of training, with a batch size of 2^{14} pixels per iteration.

We apply the sampling process outlined in Section 3.2 to align with NeRF’s hierarchical volume sampling approach [MST*20] consisting of a coarse and a fine stage. In the coarse stage, we employ stratified sampling to select points along each ray for network evaluation which involves dividing the ray into K_c equal segments and uniformly sampling a position from each segment. In the fine stage, additional samples along each ray are generated using inverse transform sampling based on density predictions from the coarse stage which results in an additional set of K_f samples.

Given a set of samples, denoted as t_1, t_2, \dots, t_K , where K corresponds to the number of coarse (K_c) or fine samples ($K_c + K_f$), TraM-NeRF establishes non-uniformly spaced intervals based on them. These intervals are defined as $\left[\frac{t_{k-1} + t_k}{2}, \frac{t_k + t_{k+1}}{2} \right]$, ensuring that each sample point t_k serves as the center of a specific section of it. We now generate directional samples by choosing a uniformly distributing length

$$t_{k,i} \sim U \left[\frac{t_{k-1} + t_k}{2}, \frac{t_k + t_{k+1}}{2} \right] \quad (18)$$

for each interval k and direction ω_i . Subsequently, the radiance field is queried for its density and color predictions at specific points

$$x_{k,i} = x + \frac{t_{k,i}}{\langle \omega^* | \omega_i \rangle} \omega_i \quad (19)$$

where the dot product between the sampled direction and the ideal reflection direction ensures that $x_{k,i}$ falls within the interval, as illustrated in Figure 3a.

4. Experiments

We ran multiple experiments to evaluate different facets of our approach, both on scenes with multiple mirrors, and scenes containing near-specular surfaces. The scenes are created in the open-source 3D graphics tool Blender to be able to accurately specify the parameters of the materials used. All 3D models and textures are provided by BlenderKit and CGTrader as royalty free assets. In total, we created 10 scenes with mirrors and 5 scenes with near-specular surfaces. For non-forward-facing scenes, we rendered 200 images per scene with cameras sampled from the upper hemisphere around the scene center, looking towards the center. For evaluation purposes, we withheld 30 per scene from the training process. As less coverage is required, for the forward-facing scenes we only rendered 100 images per scene and withheld 20 images for evaluation.

We compare our approach on our datasets, both qualitatively and quantitatively, against multiple baseline methods [MST*20; BMV*22] and recent methods that explicitly model reflections [VHM*22; GKB*22; YQCR23]. To quantify the reconstruction quality, we use the commonly used metrics peak signal-to-noise ratio (PSNR), structural similarity index (SSIM) [WBSS04] and learned perceptual image patch similarity (LPIPS) [ZIE*18].

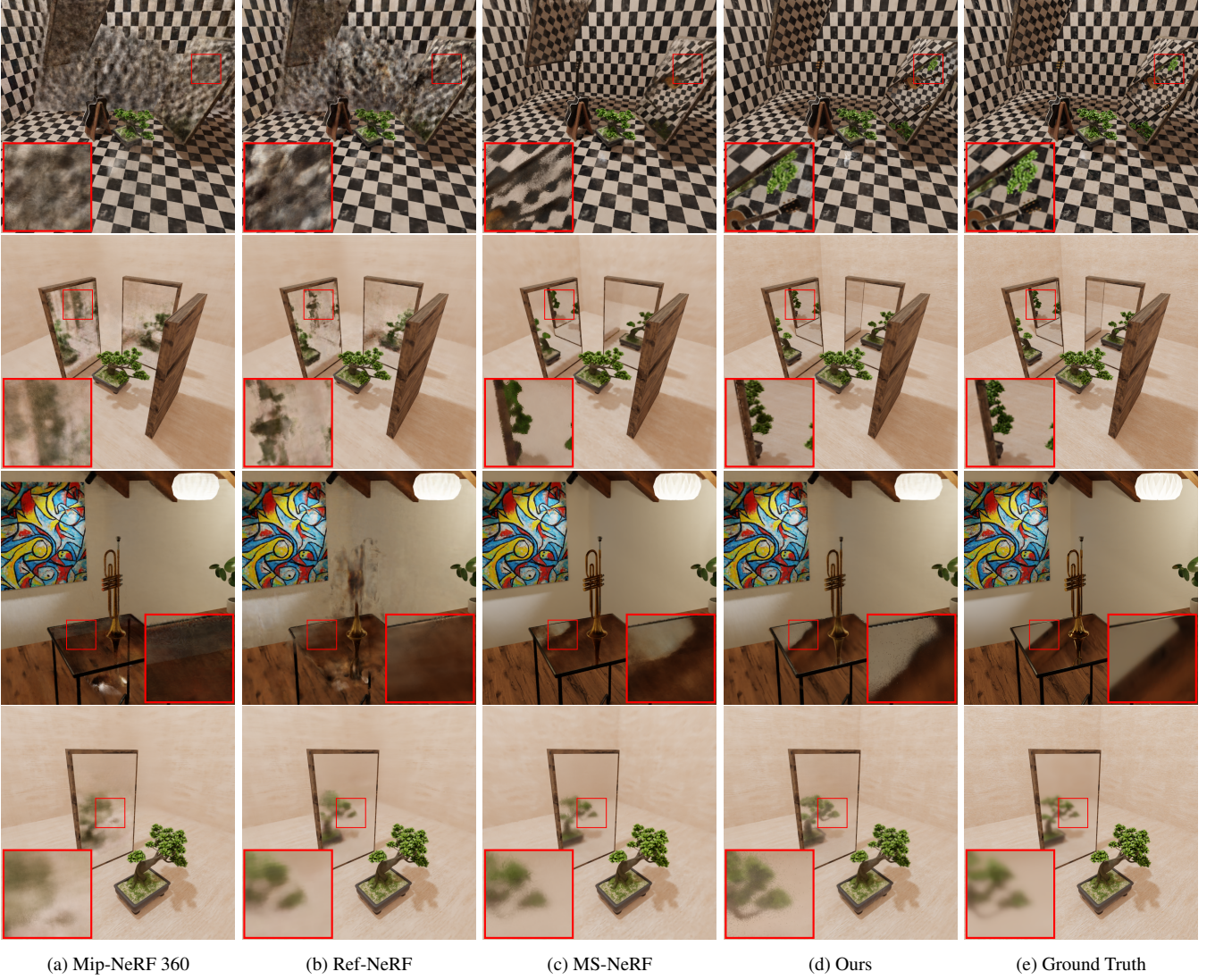


Figure 5: Comparison of different methods on our synthetic dataset with multiple mirrors (rows 1 and 2) and near-perfect specular surfaces (rows 3 and 4). The images shown are views from the test set of the respective scenes. The last column shows the ground truth test image.

4.1. Multi-Mirror Scenes

The first two rows of Figure 5 show results of various related approaches on scenes with multiple mirrors. It can be seen that both baseline methods (a) and approaches that consider reflections more explicitly (b, c) struggle to reconstruct higher-order reflections with regard to overall quality (a, b) and high-frequency details (c), while our method can by design represent these regions with the same quality as the rest of the scene. This strength is also reflected in the qualitative evaluation in Table 1, as our method outperforms the other approaches on all metrics significantly.

4.2. Reflections of Near-Specular Surfaces

The lower two rows of Figure 5 show a similar comparison on scenes with near-perfect specular surfaces. As before, (a) and (b)

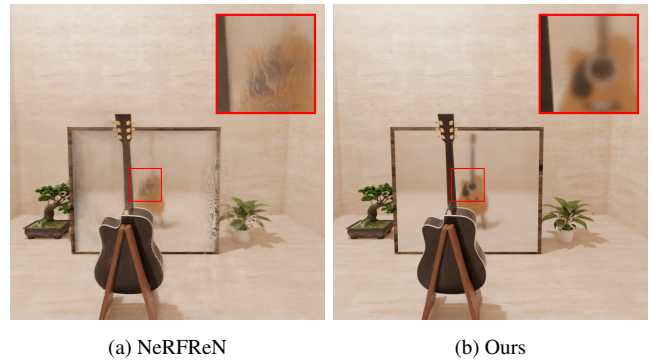


Figure 6: Results on a test view of one of the forward facing scenes that contains a near-specular surface.

	Multi-Mirror			Near-Specular Surface		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
NeRF [MST*20]	25.16 \pm 4.70	0.751 \pm 0.112	0.392 \pm 0.091	32.18 \pm 3.01	0.858 \pm 0.006	0.284 \pm 0.045
Mip-NeRF 360 [BMV*22]	24.73 \pm 5.98	0.720 \pm 0.147	0.436 \pm 0.087	32.33 \pm 2.83	0.843 \pm 0.012	0.341 \pm 0.015
Ref-NeRF [VHM*22]	24.37 \pm 5.82	0.726 \pm 0.131	0.444 \pm 0.084	31.75 \pm 4.18	0.825 \pm 0.006	0.378 \pm 0.023
MS-NeRF [YQCR23]	27.61 \pm 6.03	0.767 \pm 0.137	0.405 \pm 0.105	32.29 \pm 2.79	0.831 \pm 0.019	0.400 \pm 0.028
Ours	30.84 \pm 3.54	0.835 \pm 0.053	0.311 \pm 0.086	32.92 \pm 0.85	0.859 \pm 0.019	0.310 \pm 0.013

Table 1: Quantitative comparison of our approach against NeRF baselines and recent works on both multi-mirror scenes and scenes with near-specular surfaces. The metrics are averaged over all test images and across all scenes. Best and second-best results are highlighted.

	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
NeRF	36.33 \pm 1.54	0.904 \pm 0.017	0.212 \pm 0.021
Mip-NeRF 360	35.60 \pm 1.79	0.867 \pm 0.027	0.325 \pm 0.063
Ref-NeRF	35.51 \pm 1.68	0.860 \pm 0.028	0.344 \pm 0.058
NeRFReN	27.38 \pm 8.43	0.778 \pm 0.086	0.489 \pm 0.129
MS-NeRF	34.86 \pm 1.30	0.845 \pm 0.017	0.387 \pm 0.013
Ours	35.42 \pm 1.31	0.885 \pm 0.019	0.271 \pm 0.063

Table 2: Qualitative results on the three forward facing scenes. Best and second-best results are highlighted.

show a lack of reconstruction quality in regions close to the near-specular surface due to multi-view inconsistencies. While (c) is able to resolve the inconsistencies, in the third row it can be seen that it fails to learn a clear reflection. In both the mirror reflections and near-specular surface scenarios, we suppose that the lack of detail in the results of MS-NeRF are due to two effects: Firstly, Yin et al. reduced the sizes of the individual radiance fields to roughly match the size of approaches using a single radiance field. This leads to less capacity per radiance field in the multi-space formulation. Secondly, the previous approaches are unable to aggregate information in the reflections from different views consistently, as they are either trying to resolve the multi-view consistencies directly, or move them into a separate radiance field.

4.3. Forward-Facing Scenes

In order to compare our results with the approach of Guo et al. [GKB*22], we additionally created three scenes where all camera centers are located on a single plane. Two scenes contain mirrors, while the surface in the third scene is near-specular. The default parameters for scenes without manual annotations that are provided by the authors were used for comparison. We also experimented with providing one or multiple ground-truth masks to their approach, but we found that providing no masks consistently produced the best results.

A qualitative comparison between NeRFReN and our approach is shown in Figure 6. It can be seen that our approach is able to better reconstruct details in the near-specular region. The quantitative comparison additionally shows results of other approaches. While our approach is not reaching the highest scores on these metrics,

it closely follows the first place on the perceptual image metrics and outperforms all other methods that specifically consider reflective surfaces. This excellent performance of the baseline methods on the forward-facing scenes can be explained by the fact that this scenario does not impose multi-view inconsistencies, which are difficult to resolve using the original NeRF formulation.

4.4. Reconstruction of Indirectly Observed Regions

One of the advantages of our approach compared to works that model reflections as separate radiance fields [GKB*22; YQCR23] is that information contained in the reflection improves the reconstruction quality of the regions the reflected ray passes through. To visualize and quantify this we created additional scenes where certain regions of the scene not visible to primary camera rays in any of the training images. The cameras used to generate the test images are then chosen to cover the regions not seen in the training. An example of this setup and results are shown in Figure 7. Because the other approaches do not model a change in ray directions, they only extrapolate directly observed scene elements in the unseen regions. The periodicity in the positional encoding seems to lead to a copy of the observed scene in (b) and (d), while (c) produces noise in the respective regions. Our approach (e) on the other hand reconstructs high-frequency details that were observable in the reflection of the mirror.

4.5. Roughness Modification

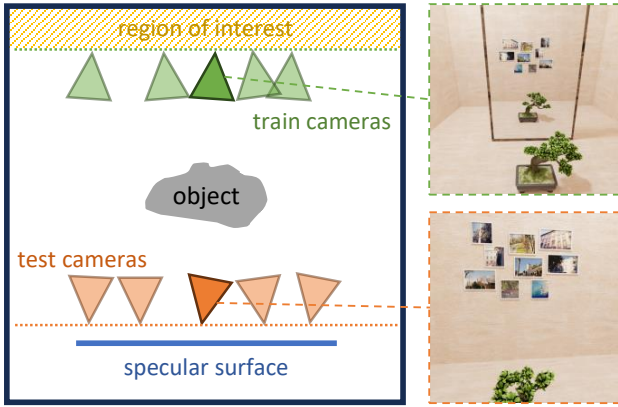
In our formulation of volumetric rendering, the BRDF parameters are decoupled from the trained radiance field, which allows us to modify these parameters at inference time. While we can increase the roughness of a mirror surface after training, we are also able to effectively de-blur a specular surface and produce a mirror-like surface at inference time. We show this change of roughness at inference time in Figure 8.

4.6. Ablation Studies

We conducted additional experiments to validate some of the design choices of our approach.

4.6.1. Annotation Robustness

To validate our claim that annotations in two images are sufficient to accurately define the position of a rectangular planar mirror in



(a) Scene Setup



(b) Mip-NeRF 360

(c) MS-NeRF



(d) Ref-NeRF

(e) Ours

Figure 7: Experiment with indirectly observed regions. (a) Schematic top view of the scene. Training cameras (green) are placed on a single plane, oriented towards a mirror (blue) that reflects light rays from an unseen region of interest (yellow) towards the training cameras. The test cameras (orange) are placed on a second plane and can directly observe the region of interest. (b)-(e) show the resulting novel views from test cameras produced by previous approaches compared to ours.

(a) $\alpha = 0$ (b) $\alpha = 0.09$ (c) $\alpha = 0.018$

Figure 8: Modification of the roughness parameter α at inference time after being trained on the value shown in (b).

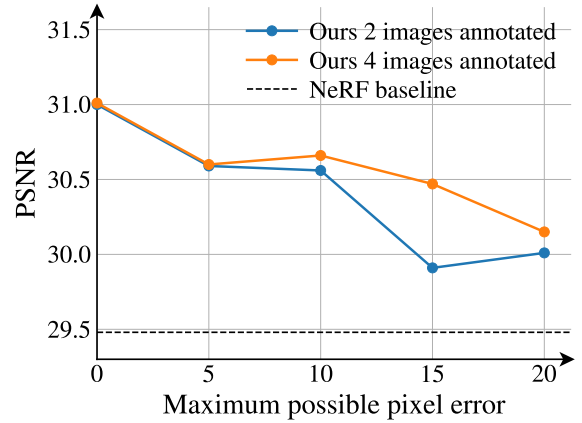


Figure 9: Mean PSNR over all test images of one of the scenes in our dataset after 10000 training iterations, using two (blue) and four (yellow) perfectly annotated images to compute the 3D location of the mirror surface. Different levels of uniform noise are added to simulate different levels of annotation errors. The dashed line shows the resulting mean PSNR of NeRF for the exact annotations after the same amount of training iterations.

3D space, we perturbed the annotated positions in screen space with noise sampled from a uniform distribution over the discrete interval $[-k, k]$. We varied the parameter k to simulate increasingly severe levels of errors in the annotation. Figure 9 shows the resulting PSNR of the reconstruction with pixel-exact annotations and different levels of noise on a single scene in our dataset. It can be seen that carefully annotating only two images yield comparable results to adding two more annotated images. While higher noise levels in the annotation leads to a drop in reconstruction quality, our approach shows to be robust as it is able to yield an improvement compared to the NeRF baseline with pixel-exact annotations, even under severe noise.

4.6.2. Modified Ray Sampling

As motivated in Section 3.2, we draw the microfacet normals independently for each sampling location along the main ray (dense) instead of generating multiple rays at the surface intersection and choosing sampling points along these rays (sparse). We ran an experiment to compare the reconstruction quality of these variants. Figure 10 shows that the sparse variant introduces significantly

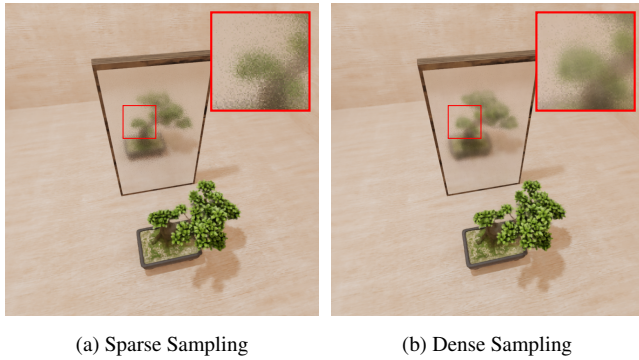


Figure 10: Comparison of sparse and dense sampling variants of our approach on one of the scene of our dataset.

more noise on the near-specular surface than the modified, dense variant. This, in turn, shows that more samples would be required when using the sparse sampling method to achieve a similar level of noise.

4.7. Limitations

While TraM-NeRF achieves promising results for novel view synthesis, it also has some limitations that require further attention. In the context of generating novel views, our model inherits NeRF’s limitations in extrapolating effectively in areas with insufficient input image coverage, leading to reduced performance. We observed these hallucinations in parts of the scene that are concealed behind mirrors in a majority of training images. Additionally, our estimator can overestimate density and transmittance when the assumption of a narrow spread of light directions does not hold. In this case object reflections can appear larger than expected. Moreover, the current estimator implementation is limited to single mirror-like surfaces.

5. Conclusions

We presented TraM-NeRF, an extension of NeRF that effectively models mirror-like surfaces, accurately capturing high-frequency reflections within a single scene representation. By introducing a transmittance-aware variant of the rendering equation for explicit reflection modeling as well as efficient sampling techniques, we are able to reduce the number of network evaluations during ray tracing without increasing the variance. In the scope of a qualitative and quantitative evaluation, we demonstrated that our techniques outperform previous methods in challenging scenes with single and multiple mirror-like surfaces.

Acknowledgements

This work has been funded by the DFG project KL 1142/11-2 (DFG Research Unit FOR 2535 Anticipating Human Behavior), and additionally by the Federal Ministry of Education and Research of Germany and the state of North-Rhine Westphalia as part of the Lamarr-Institute for Machine Learning and Artificial Intelligence and by the Federal Ministry of Education and Research under grant no. 01IS22094E WEST-AI.

References

- [BBJ*21] BOSS, M., BRAUN, R., JAMPANI, V., et al. “NeRD: Neural Reflectance Decomposition from Image Collections”. *IEEE International Conference on Computer Vision (ICCV)*. 2021, 12684–12694 2.
- [BEK*22] BOSS, M., ENGELHARDT, A., KAR, A., et al. “Samurai: Shape and material from unconstrained real-world arbitrary image collections”. *Advances in Neural Information Processing Systems (NeurIPS)* 35 (2022), 26389–26403 2.
- [BMT*21] BARRON, J. T., MILDENHALL, B., TANCİK, M., et al. “Mip-NeRF: A Multiscale Representation for Anti-Aliasing Neural Radiance Fields”. *IEEE International Conference on Computer Vision (ICCV)*. 2021, 5855–5864 1, 2, 6.
- [BMV*22] BARRON, J. T., MILDENHALL, B., VERBIN, D., et al. “Mip-NeRF 360: Unbounded Anti-Aliased Neural Radiance Fields”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2022, 5470–5479 2, 6, 8.
- [BXS*20a] BI, S., XU, Z., SRINIVASAN, P., et al. “Neural Reflectance Fields for Appearance Acquisition”. *arXiv preprint arXiv:2008.03824* (2020) 2.
- [BXS*20b] BI, S., XU, Z., SUNKAVALLI, K., et al. “Deep reflectance volumes: Relightable reconstructions from multi-view photometric images”. *European Conference on Computer Vision (ECCV)*. Springer. 2020, 294–311 2.
- [CT82] COOK, R. L. and TORRANCE, K. E. “A reflectance model for computer graphics”. *ACM Transactions on Graphics (ToG)* 1.1 (1982), 7–24 5.
- [CXG*22] CHEN, A., XU, Z., GEIGER, A., et al. “Tensorf: Tensorial radiance fields”. *European Conference on Computer Vision (ECCV)*. 2022, 333–350 2, 6.
- [CZL*22] CHEN, X., ZHANG, Q., LI, X., et al. “Hallucinated neural radiance fields in the wild”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2022, 12943–12952 1, 2.
- [DB23] DUPUY, J. and BENYOUB, A. “Sampling Visible GGX Normals with Spherical Caps”. *Computer Graphics Forum (CGF)* (2023) 5.
- [FMW*23] FRIDOVICH-KEIL, S., MEANTI, G., WARBURG, F. R., et al. “K-planes: Explicit radiance fields in space, time, and appearance”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2023, 12479–12488 2.
- [FSV*23] FAN, Y., SKOROKHODOV, I., VOYNOV, O., et al. “Factored-NeuS: Reconstructing Surfaces, Illumination, and Materials of Possibly Glossy Objects”. *arXiv preprint arXiv:2305.17929* (2023) 2.
- [FYT*22] FRIDOVICH-KEIL, S., YU, A., TANCİK, M., et al. “Plenoxels: Radiance Fields Without Neural Networks”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2022, 5501–5510 2.
- [GHZ*23] GE, W., HU, T., ZHAO, H., et al. “Ref-NeuS: Ambiguity-Reduced Neural Implicit Surface Learning for Multi-View Reconstruction with Reflection”. (2023) 2.
- [GKB*22] GUO, Y.-C., KANG, D., BAO, L., et al. “NeRFReN: Neural Radiance Fields With Reflections”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2022, 18409–18418 2, 3, 6, 8.
- [GKJ*21] GARBIN, S. J., KOWALSKI, M., JOHNSON, M., et al. “Fastnerf: High-fidelity neural rendering at 200fps”. *IEEE International Conference on Computer Vision (ICCV)*. 2021, 14346–14355 2.
- [Hei18] HEITZ, E. “Sampling the GGX distribution of visible normals”. *Journal of Computer Graphics Techniques (JCGT)* 7.4 (2018), 1–13 5.
- [HZF*22] HUANG, X., ZHANG, Q., FENG, Y., et al. “Hdr-nerf: High dynamic range neural radiance fields”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2022, 18398–18408 2.
- [JLX*23] JIN, H., LIU, I., XU, P., et al. “TensorIR: Tensorial Inverse Rendering”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2023, 165–174 2.

- [Kaj86] KAJIYA, J. T. “The rendering equation”. *Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*. 1986, 143–150 [4](#).
- [LCL*23] LIANG, R., CHEN, H., LI, C., et al. “ENVIDR: Implicit Differentiable Renderer with Neural Environment Lighting”. *IEEE International Conference on Computer Vision (ICCV)*. 2023 [2](#).
- [LLGG23] LI, Q., LI, F., GUO, J., and GUO, Y. “UHDNeRF: Ultra-High-Definition Neural Radiance Fields”. *IEEE International Conference on Computer Vision (ICCV)*. 2023, 23097–23108 [2](#).
- [LSS*19] LOMBARDI, S., SIMON, T., SARAGIH, J., et al. “Neural volumes: learning dynamic renderable volumes from images”. *ACM Transactions on Graphics (TOG)* 38.4 (2019), 1–14 [2](#).
- [Max95] MAX, N. “Optical models for direct volume rendering”. *IEEE Transactions on Visualization and Computer Graphics (TVCG)* 1.2 (1995), 99–108 [3](#).
- [MESK22] MÜLLER, T., EVANS, A., SCHIED, C., and KELLER, A. “Instant Neural Graphics Primitives with a Multiresolution Hash Encoding”. *ACM Transactions on Graphics (TOG)* 41.4 (2022), 102:1–102:15 [1](#), [2](#), [6](#).
- [MHM*22] MILDENHALL, B., HEDMAN, P., MARTIN-BRUALLA, R., et al. “Nerf in the dark: High dynamic range view synthesis from noisy raw images”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2022, 16190–16199 [2](#).
- [MRS*21] MARTIN-BRUALLA, R., RADWAN, N., SAJJADI, M. S., et al. “Nerf in the wild: Neural radiance fields for unconstrained photo collections”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2021, 7210–7219 [2](#).
- [MST*20] MILDENHALL, B., SRINIVASAN, P. P., TANCIK, M., et al. “NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis”. *European Conference on Computer Vision (ECCV)*. 2020 [1](#)–[6](#), [8](#).
- [MT97] MÖLLER, T. and TRUMBORE, B. “Fast, Minimum Storage Ray-Triangle Intersection”. *Journal of Graphics Tools* 2.1 (1997), 21–28 [6](#).
- [MVKF23] MAI, A., VERBIN, D., KUESTER, F., and FRIDOVICH-KEIL, S. “Neural Microfacet Fields for Inverse Rendering”. *IEEE International Conference on Computer Vision (ICCV)*. 2023, 408–418 [2](#), [3](#).
- [NMOG20] NIEMEYER, M., MESCHEDER, L., OECHSLE, M., and GEIGER, A. “Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2020, 3504–3515 [2](#).
- [PGM*19] PASZKE, A., GROSS, S., MASSA, F., et al. “PyTorch: An Imperative Style, High-Performance Deep Learning Library”. *Advances in Neural Information Processing Systems (NeurIPS)*. 2019, 8024–8035 [6](#).
- [RPLG21] REISER, C., PENG, S., LIAO, Y., and GEIGER, A. “Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps”. *IEEE International Conference on Computer Vision (ICCV)*. 2021, 14335–14345 [1](#), [2](#).
- [SDZ*21] SRINIVASAN, P. P., DENG, B., ZHANG, X., et al. “NeRV: Neural Reflectance and Visibility Fields for Relighting and View Synthesis”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2021, 7495–7504 [2](#).
- [SZW19] SITZMANN, V., ZOLLHÖFER, M., and WETZSTEIN, G. “Scene representation networks: Continuous 3d-structure-aware neural scene representations”. *Advances in Neural Information Processing Systems (NeurIPS)* 32 (2019) [2](#).
- [VHM*22] VERBIN, D., HEDMAN, P., MILDENHALL, B., et al. “Ref-NeRF: Structured View-Dependent Appearance for Neural Radiance Fields”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE. 2022, 5481–5490 [2](#), [6](#), [8](#).
- [WBSS04] WANG, Z., BOVIK, A. C., SHEIKH, H. R., and SIMONCELLI, E. P. “Image quality assessment: from error visibility to structural similarity”. *IEEE Transactions on Image Processing (TIP)* 13.4 (2004), 600–612 [6](#).
- [WHL*23] WU, H., HU, Z., LI, L., et al. “NeFII: Inverse Rendering for Reflectance Decomposition with Near-Field Indirect Illumination”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2023, 4295–4304 [2](#).
- [WLS*22] WANG, Z., LI, L., SHEN, Z., et al. “4K-NeRF: High fidelity neural radiance fields at ultra high resolutions”. *arXiv preprint arXiv:2212.04701* (2022) [2](#).
- [WMLT07] WALTER, B., MARSCHNER, S. R., LI, H., and TORRANCE, K. E. “Microfacet models for refraction through rough surfaces”. *Eurographics Symposium on Rendering (EGSR)*. 2007, 195–206 [5](#).
- [WWG*22] WANG, C., WU, X., GUO, Y.-C., et al. “NeRF-SR: High Quality Neural Radiance Fields using Supersampling”. *ACM International Conference on Multimedia*. 2022, 6445–6454 [1](#), [2](#).
- [YQCR23] YIN, Z.-X., QIU, J., CHENG, M.-M., and REN, B. “Multi-Space Neural Radiance Fields”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2023, 12407–12416 [2](#), [3](#), [6](#), [8](#).
- [ZBC*23] ZENG, J., BAO, C., CHEN, R., et al. “Mirror-NeRF: Learning Neural Radiance Fields for Mirrors with Whitted-Style Ray Tracing”. *ACM International Conference on Multimedia*. 2023 [2](#), [3](#).
- [ZIE*18] ZHANG, R., ISOLA, P., EFROS, A. A., et al. “The unreasonable effectiveness of deep features as a perceptual metric”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2018, 586–595 [6](#).
- [ZLW*21] ZHANG, K., LUAN, F., WANG, Q., et al. “PhysSG: Inverse Rendering with Spherical Gaussians for Physics-based Material Editing and Relighting”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2021, 5453–5462 [2](#).
- [ZRSK20] ZHANG, K., RIEGLER, G., SNAVELY, N., and KOLTUN, V. “Nerf++: Analyzing and improving neural radiance fields”. *arXiv preprint arXiv:2010.07492* (2020) [1](#), [2](#).
- [ZSD*21] ZHANG, X., SRINIVASAN, P. P., DENG, B., et al. “NeRFactor: Neural Factorization of Shape and Reflectance Under an Unknown Illumination”. *ACM Transactions on Graphics (TOG)* 40.6 (2021), 1–18 [2](#).
- [ZSH*22] ZHANG, Y., SUN, J., HE, X., et al. “Modeling indirect illumination for inverse rendering”. *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2022, 18643–18652 [2](#).
- [ZXY*23] ZHANG, Y., XU, T., YU, J., et al. “NeMF: Inverse Volume Rendering with Neural Microflake Field”. *IEEE International Conference on Computer Vision (ICCV)*. 2023, 22919–22929 [2](#), [3](#).