

FG-NeRF: Flow-GAN based Probabilistic Neural Radiance Field for Independence-Assumption-Free Uncertainty Estimation

Songlin Wei^{*1,2}, Jiazhao Zhang^{*1,2}, Yang Wang¹, Fanbo Xiang³, Hao Su³, He Wang^{1,2}

¹CFCS, School of Computer Science, Peking University

²Beijing Academy of Artificial Intelligence

³UC San Diego

Abstract

Neural radiance fields with stochasticity have garnered significant interest by enabling the sampling of plausible radiance fields and quantifying uncertainty for downstream tasks. Existing works rely on the independence assumption of points in the radiance field or the pixels in input views to obtain tractable forms of the probability density function. However, this assumption inadvertently impacts performance when dealing with intricate geometry and texture. In this work, we propose an independence-assumption-free probabilistic neural radiance field based on Flow-GAN. By combining the generative capability of adversarial learning and the powerful expressivity of normalizing flow, our method explicitly models the density-radiance distribution of the whole scene. We represent our probabilistic NeRF as a mean-shifted probabilistic residual neural model. Our model is trained without an explicit likelihood function, thereby avoiding the independence assumption. Specifically, We downsample the training images with different strides and centers to form fixed-size patches which are used to train the generator with patch-based adversarial learning. Through extensive experiments, our method demonstrates state-of-the-art performance by predicting lower rendering errors and more reliable uncertainty on both synthetic and real-world datasets.

Introduction

As a popular approach for neural scene modeling, Neural Radiance Fields (NeRF) (Mildenhall et al. 2021) have been widely studied in recent years due to their ability to render photorealistic novel views. Additionally, NeRF has shown effectiveness in various related fields, including 3D reconstruction (Wang et al. 2021; Yariv et al. 2021; Müller et al. 2022), camera-pose recovery (Meng et al. 2021; Schwarz et al. 2020), 3D semantic segmentation (Zhi et al. 2021; Vora et al. 2021), and photo-realistic simulation (Li et al. 2023, 2021). However, the majority of existing works (Gao et al. 2022) only output one color for each pixel while falling short to capture the uncertainty in neural scene modeling. This limitation hinders the applications of NeRFs in active perception and interaction tasks, *e.g.*, robotics (Ran et al. 2022; Pan et al. 2022), autonomous driving (Li et al. 2022b), and human-computer interaction (Li et al. 2022a), where it is essential to incorporate uncertain information (Cai et al. 2021).

^{*}These authors contributed equally.

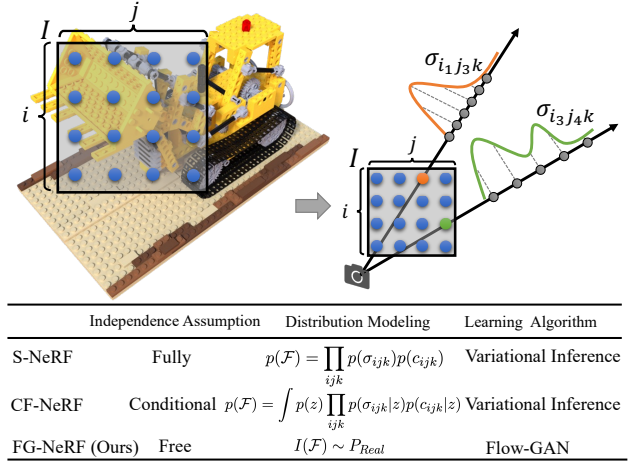


Figure 1: We propose to use adversarial learning to model probabilistic neural radiance field \mathcal{F} . Previous work models the whole density-radiance probability distribution as $p(\mathcal{F})$. S-NeRF factorizes the $p(\mathcal{F})$ as the product of the density distribution $p(\sigma_{ijk})$ and radiance distribution $p(c_{ijk})$ of all sampled points, while CF-NeRF (Shen et al. 2022) uses latent variable z to make the sampled points conditional independent. On the contrary, our work FG-NeRF directly models the distribution P_{Real} of the rendered image I through volume rendering of \mathcal{F} .

To capture the uncertainty of NeRF, recent studies have proposed various types of probabilistic neural radiance fields (Shen et al. 2021, 2022) that aim to learn a distribution over the radiance field. However, learning a distribution over such continuous fields presents significant challenges, due to the vast probability space involved. Consequently, existing approaches usually rely on strong assumptions. For example, Stochastic NeRF (S-NeRF) (Shen et al. 2021) assumes the radiance and density values are independent Gaussian distributions at different locations, which neglects the important spatial consistency in 3D spaces. The state-of-the-art method, Conditional-Flow NeRF (CF-NeRF) (Shen et al. 2022), leverages normalizing flow (Winkler et al. 2019) technique to capture arbitrary radiance/density distributions. Furthermore, CF-NeRF introduces a global latent variable

and achieves conditional independence among the values via De Finetti’s representation theorem (Kirsch 2019). However, we contend that the conditional independence assumption still restricts the capacity of probabilistic NeRF to represent complex and large scenes accurately.

In this work, we propose a novel probabilistic neural radiance field, Flow-GAN NeRF (FG-NeRF). Inspired by Flow-GAN (Grover, Dhar, and Ermon 2018a), our method leverages both normalizing flow (Winkler et al. 2019) and GAN (Goodfellow et al. 2020). We propose to use adversarial training to learn the radiance/density distributions of the whole scene. In our FG-NeRF, the normalizing flow module serves only as a generator and is trained along with an additional discriminator, which together form a GAN framework. As a result, our method exhibits two primary distinctions from S-NeRF or CF-NeRF (See Fig. 1: *first*, we leverage generative adversarial learning instead of variational inference used in prior methods. *Second*, our method dispenses with the need for scene-wide distribution factorization to compute the likelihood or self-entropy term during variational inference (See Sec. 5), thereby achieving independence-assumption free. Consequently, FG-NeRF demonstrates the capability to produce high-quality uncertainty estimations while effectively capturing intricate appearance and geometry in room-level scenes.

Through extensive experiments, FG-NeRF showcases new state-of-the-art performance when compared to existing uncertainty estimation methods on LLFF (Mildenhall et al. 2021), ScanNet (Dai et al. 2017), and Replica (Straub et al. 2019) datasets. Our proposed method achieves improved uncertainty estimation, accompanied by high-quality color and depth predictions. Moreover, by integrating FG-NeRF with a multi-level hash encoding technique (Müller et al. 2022), our approach achieves a remarkable four times faster compared to CF-NeRF. With this significant improvement in computational efficiency, we further build an autonomous NeRF reconstruction, showcasing the superior performance of our predicted uncertainty.

Related Work

NeRF (Mildenhall et al. 2021) has been widely studied due to its impressive capability in reconstructing scenes for photorealistic novel views. By utilizing a sparse collection of 2D views and employing a neural network to encode a representation of the 3D scene. This neural network consists of a set of fully-connected layers that output the radiance for any given input 3D spatial location and 2D viewing direction within the scene. In recent years, NeRF has been extensively researched to improve the quality (Zhang et al. 2020; Kosiorek et al. 2021; Barron et al. 2021; Wang et al. 2021) and efficiency (Neff et al. 2021; Müller et al. 2022; Chen et al. 2022a). However, despite the advancements made, NeRF is unable to quantify the uncertainty associated with rendered views or estimated geometry, which limits its applicability in active tasks. To address this limitation, we propose FG-NeRF, a probabilistic neural radiance field that explicitly quantifies uncertainty.

Uncertainty estimation in NeRF. The estimation of uncertainty has been a widely explored topic in deep learning (Hendrycks et al. 2019; Brachmann et al. 2016; Malinin and Gales 2018; Teye, Azizpour, and Smith 2018; Geifman, Uziel, and El-Yaniv 2018). Here, we only review the papers that are highly relevant to NeRF. Early studies (Lee et al. 2022) regard uncertainty as the entropy of the density among the sampled ray. Despite the simplicity, this approach tends to perform poorly on complex geometry and lacks the capability to estimate uncertainty at arbitrary positions. Another line of research (Ran et al. 2022; Pan et al. 2022; Martin-Brualla et al. 2021) attempted to implicitly predict the variance associated with rendered pixels or positions. However, These methods apply the volume rendering process over uncertainty values are not theoretically grounded.

Recently, the most compelling approach is to learn a posterior distribution over the model parameters. S-NeRF (Shen et al. 2021), for instance, models a distribution encompassing all possible radiance fields that explain the scene. To ensure tractability, S-NeRF makes a strong assumption by assuming that all predicted distributions of each point are independent. In an effort to alleviate this limitation, CF-NeRF (Shen et al. 2022) softens this assumption by conditioning the predicted distributions on a shared latent variable (Yang et al. 2019). This design enables spatially smooth distribution and improves the quality of predictions. However, we argue that CF-NeRF still relies on a conditional independence assumption through the use of a global latent variable. Notably, our FG-NeRF stands out by sidestepping the independence assumptions by employing a generator and a discriminator that are trained adversarially.

Adversarial Learning for NeRF. Adversarial learning (Goodfellow et al. 2020; Creswell et al. 2018; Grover, Dhar, and Ermon 2018b) has recently emerged as a promising approach for learning the radiance field for various challenging tasks. Previous relevant studies (Meng et al. 2021; Schwarz et al. 2020; Chan et al. 2021) have primarily focused on leveraging generative models to learn the radiance field from unstructured 2D images. Recent research efforts (Deng et al. 2022; Chan et al. 2022) have made significant advancements in enhancing synthesis quality by achieving higher resolutions, improved 3D geometry, and faster rendering. However, to the best of our knowledge, our work is the first to utilize adversarial learning specifically for uncertainty estimation in NeRF.

Background

Deterministic and probabilistic formulations of NeRF

Neural radiance field (NeRF) (Mildenhall et al. 2021) is an implicit scene representation. It is modeled as a continuous function $\mathcal{F} : (x, d) \rightarrow (\sigma, c)$, which maps the pair of 3D spatial location $x \in \mathbb{R}^3$ and viewing direction $d \in \mathbb{S}^2$ into density $\sigma \in \mathbb{R}^+$ and RGB radiance $c \in [0, 1]^3$. Given a camera ray $r(t) = x_0 + td$ shooting from origin x_0 along the direction d , the color c of the ray is computed via the

physics-based volume rendering:

$$c(r) = \int_0^\infty T(t) \sigma(r(t)) c(r(t), d) dt, \quad (1)$$

where $T(t) = \exp(-\int_0^t \sigma(r(s)) ds)$,

and the function $T(t)$ is the accumulated transmittance along the ray. It is the probability that the ray travels from the origin to x_t without hitting any obstacles.

The radiance field \mathcal{F} is *deterministic* in terms of having a single density and RGB radiance for each point in the space. Deterministic NeRF is unable to output the uncertainty of the rendering when the output confidence is desired. One approach to remedy this problem is to formulate the neural radiance construction as *probabilistic*. Probabilistic NeRF models the scene as a joint probabilistic density and radiance distribution $p(\mathcal{F})$ for all the spatial locations and view directions. In this manner, the deterministic radiance field can be regarded as one realization over the distribution $p(\mathcal{F})$.

Modeling density-radiance distributions with normalizing flow To learn the probability distribution \mathcal{F} of the whole scene, we first discuss how to model the distribution of a single point through normalizing flow. For any location x viewed from direction d , the conditional joint distribution of the density and radiance is $p(\sigma, c|x, d) = p(\sigma|x)p(c|x, d)$ where density σ is assumed to depend only on location x . Approximating the radiance distribution $p(c|x, d)$ with unimodal distributions like Gaussian is sub-optimal, as a point in the real-world can have different colors viewed from different directions due to reflection or iridescence. Inspired by CF-NeRF (Shen et al. 2022), we use conditional normalizing flow (CNF) (Winkler et al. 2019) which transforms a simple distribution into arbitrary complex shape conditioning on the location-view pair (x, d) . Specifically, conditional normalizing flow maps the initial random variable (σ_0, c_0) with the prior distributions $q_\sigma(\sigma_0)$ and $q_c(c_0)$ into the target (σ, c) via K steps of bijective mappings $\sigma_k = f_k^\sigma(\sigma_{k-1}, x)$ and $c_k = f_k^c(c_{k-1}, x, d)$ respectively. The target distributions can be computed via the change-of-variables formula:

$$p(\sigma|x) = q_\sigma(\sigma_0) \prod_{k=1}^{K-1} \left| \det \frac{\partial f_k^\sigma}{\partial \sigma_{k-1}} \right|^{-1}, \quad (2)$$

$$p(c|x, d) = q_c(c_0) \prod_{k=1}^{K-1} \left| \det \frac{\partial f_k^c}{\partial c_{k-1}} \right|^{-1} \quad (3)$$

The initial distribution q_σ and q_c is usually set to standard Gaussian distributions.

Because directly computing the density-radiance distribution $p(\mathcal{F})$ for the whole scene is intractable, previous work (Shen et al. 2021)(Shen et al. 2022) approximate the probabilistic distribution using variational inference (Blei, Kucukelbir, and McAuliffe 2017) via a parametric function $q(\mathcal{F}; \theta)$. The models boil down to the computation of a likelihood term $p(\mathcal{T}|\mathcal{F})$ of the training sample and a self-entropy term. They assume the distribution for each point in the scene is independent or conditionally independent to

constrain the $q(\mathcal{F}; \theta)$, which can then be factorized into the product of each point distribution to make the computation tractable:

$$q_\theta(\mathcal{F}) = \prod_{x \in \mathbb{R}^3} \prod_{d \in \mathbb{S}^2} p(\sigma|x)p(c|x, d), \quad \text{or} \quad (4)$$

$$q_\theta(\mathcal{F}|z) = \int p(z) \prod_{x \in \mathbb{R}^3} \prod_{d \in \mathbb{S}^2} p(\sigma|x, z)p(c|x, d, z). \quad (5)$$

Nevertheless, neighboring 3D points within the scene are prone to exhibit analogous density or color, thereby casting doubt on the plausibility of the conditional independence assumption.

Although these works produce high-quality results of uncertainty estimation on object-level scenes, experiments show that this assumption hampers the performance when the scene becomes large and complex.

Method

We study the probability distribution of the entire density-radiance field. More concretely, given $\mathcal{T} = \{(I^i, T^i, D^i)\}_{i=1}$ denotes the triplets of image I , pose T and optional depth D indexed by i , we goal is to learn the posterior distribution $p(\mathcal{F}|\mathcal{T})$.

Our whole pipeline is illustrated in Fig. 2. Given an image I , we dynamically sample a batch of rays \mathcal{R} to obtain an image patch P . The location-view pairs (x, d) which are sampled through the stratified sampling of the rays originating from camera pose T are fed into the generator G_θ . The generator consists of two branches: (1) the deterministic branch learns the mean values for (2) the probabilistic branch which, is implemented with the conditional normalizing flow (CNF), learns to predict the mean-shifted probability distribution of the radiance field of the scene. Finally, a synthesized patch \hat{P} is obtained through volume rendering of the combined branches and is sent to the discriminator D_ϕ to compare with the real training patch P .

Decomposing probabilistic neural radiance field

Directly learning the probability distribution of the density-radiance field presents additional challenges due to the significant divergences in distribution between occupied and free regions. Consequently, the presence of large density values (ranging from 0 to $+\infty$) can lead to numerical instability during the training of the normalizing flow. To tackle this challenge, we decompose \mathcal{F} into a deterministic branch and a residual probabilistic branch:

$$\mathcal{F} = \mathcal{F}^\mu + \mathcal{F}^\epsilon \quad (6)$$

The deterministic branch \mathcal{F}^μ acting like the mean of the probabilistic density-radiance field is estimated following the standard volume rendering procedure. It provides the mean values of density and radiance for each location-view pair (x, d) in the scene. Thus, the probabilistic branch \mathcal{F}^ϵ is only needed to learn a mean-shifted residual distribution. The distributions are learned via conditional normalizing flow as described in Section . Given a batch of rays \mathcal{R} sampled from the training images I , the deterministic branch

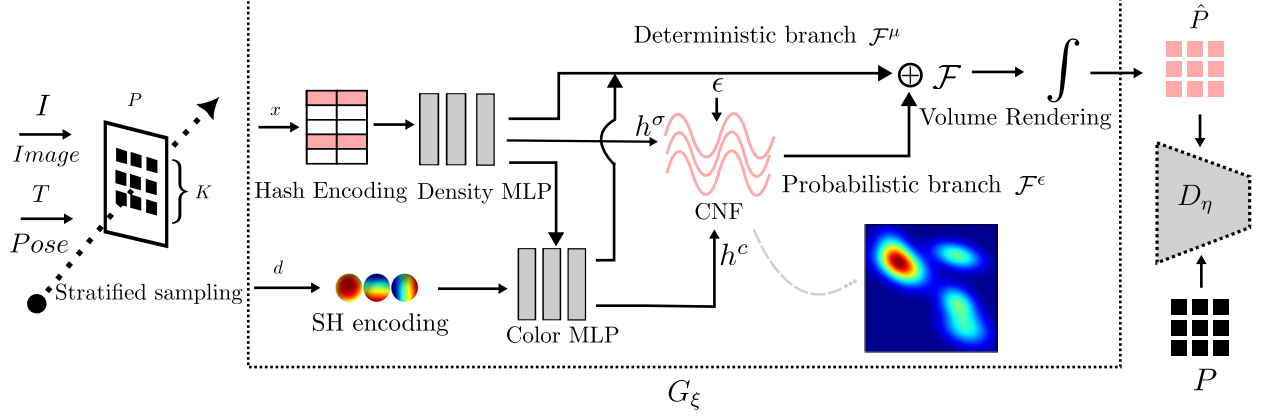


Figure 2: Overview of our pipeline. Our pipeline learns to generate fake image patches given the training image I and pose T . The generator G_ξ and discriminator D_η jointly guide the normalizing flow (CNF) for modeling the density/radiance distribution.

\mathcal{F}^μ of the network is supervised with the $L2$ loss between the rendered color $\hat{c}(r)$ and the ground truth color of each ray r . The loss is denoted as

$$\mathcal{L}_{det} = \sum_{r \in \mathcal{R}} \|I[r] - \hat{c}(r)\|_2^2 \quad (7)$$

where $\hat{c}(r)$ represents the rendered color of the ray r . It is obtained through alpha-composition:

$$\begin{aligned} \hat{c} &= \sum_{i=1}^N T^i \alpha^i c^{\mu, i}, \\ T^i &= \prod_{j=1}^{i-1} (1 - \alpha^j), \\ \alpha^i &= 1 - \exp(-\sigma^{\mu, i} \delta^i) \end{aligned} \quad (8)$$

where superscript i indicates the index of the discrete sampling points along the ray. α denotes the opacity of the color at the sampling point and δ is the distance between neighboring sampling points. With this deterministic branch been regularized by the rendering loss, our network can efficiently learn a rough geometry and appearance of the scene. Furthermore, the residual mean-shifted probabilistic branch \mathcal{F}^ϵ is eased to learn a zero-mean centered distribution of density σ^ϵ and radiance c^ϵ .

Thus, the combined density and radiance are

$$\sigma = \sigma^\mu + \sigma^\epsilon \quad c = c^\mu + c^\epsilon. \quad (9)$$

Adversarial learning of probabilistic NeRF based on Flow-GAN

Instead of the factorizations of the probability distributions as done in (Shen et al. 2021)(Shen et al. 2022), we propose to employ adversarial learning to sidestep the direct computation of the likelihood functions as illustrated in Fig. 3. Therefore, our method does not rely on any form of independence assumptions. We first dynamically sample pixels (Isola et al. 2017) structured with $K \times K$ pixels from training image I . The image patch follows the real data distribution p_{data} . Then, we train a generator G_ξ to synthesize the

patch using known camera pose T and a discriminator D_η to compare the real and fake patches. More specifically, the

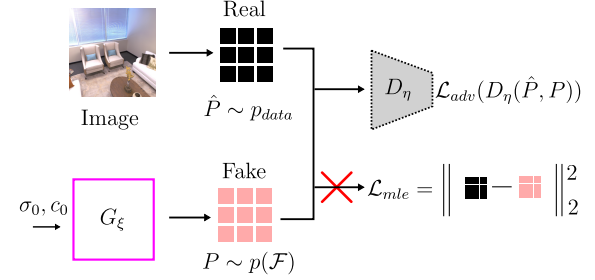


Figure 3: Comparison of patch-based adversarial loss and pixel-level rendering loss

K^2 rays in the patch are spaced evenly with random intervals and the center is also randomly shifted within the image. This sampling strategy preserves the image structure which allows the generator to generate high-quality samples. We then employ stratified sampling following standard NeRF to get N points along each ray. We follow (Müller et al. 2022) to use multi-resolution hash encoding to encode sampling points x and use the spherical harmonics to encode the view direction d .

The generator G_ξ is implemented with CNF. For each sampling point x , we first acquire its mean density σ^μ and color c^μ through the deterministic branch \mathcal{F}^μ and then transform an initial density σ^ϵ and color c^ϵ following standard Gaussians into the target distribution using CNF, which is conditional on the location feature h^σ and radiance feature h^c . Subsequently, we synthesize the image patch via volume rendering of the combined densities and colors as in Eq. (9) and Eq. (8). The synthesized image \hat{P} will be fed into a discriminator D_η to be compared with the real image patches P . Finally, we minimize a distribution distance between the generated image $p(\mathcal{F})$ and the real image patch p_{data} based on adversarial loss \mathcal{L}_{adv} as is done in classic GAN frame-

work. The training objective is:

$$\min_{\xi \in \Xi} \max_{\eta \in E} \mathcal{L}_{adv} = \mathbb{E}_{P \sim p_{data}} [\log(D(P; \eta))] + \mathbb{E}_{\hat{P} \sim p(\mathcal{F})} [\log(1 - D(G(\sigma_0, c_0; \xi); \eta))]. \quad (10)$$

The discriminator D_η is trained to output 1 for real patches and 0 for fake patches. It is implemented as a convolutional neural network following (Schwarz et al. 2020). The discriminator is conditioned on the scale s to perform well. The scale s is the ratio of the size of the sampled patch with respect to the input image. s is set to 1 at the beginning of the training to have a large perceptive field and gradually decreased to a small value to capture fine-grained details of the scene. The generator G_ξ and discriminator D_η are trained adversarially to promote the generator to synthesize more realistic images. As a result, the normalizing flow in the generator is trained to produce high-quality density-radiance distributions for synthesizing photo-realistic images. Thus, uncertainty quantified by such distributions has better qualities compared with previous works as demonstrated in the experiments.

Optional extra depth supervision

Unlike object-level scenes where training rays are densely available for every point in the scene. It is challenging for NeRF to model large scenes like ScanNet (Dai et al. 2017) due to the limited number of training images. In such cases, depth supervision further improves the learning of geometry and radiance. The ray depth can be obtained by the integral:

$$\hat{D} = \sum_{i=1}^N w^i R^i \quad (11)$$

where $w^i = T^i \alpha^i$ is the weight of the ray depth R^i of the sampling point. However, straightforward $L2$ loss $\mathcal{L}_{depth} = \|D - \hat{D}\|_2^2$ supervision performs poorly. We empirically found that using the cross-entropy loss between the sample weights $w(t) = T(t)\alpha(t)$ and an assumed ground truth depth centered Gaussian distribution $\mathcal{N}(D, \sigma^2)$ performs better than straightforward $L2$ loss. The cross-entropy loss is defined as:

$$\mathcal{L}_{depth} = \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \log w_r^i \exp\left(-\frac{(t_i - D)^2}{2\sigma^2}\right) \quad (12)$$

where Σ is set to 0.1 in the experiments. Intuitively, such cross entropy constrains the weights to follow a unimodal distribution hence enforcing stronger regularization when learning the density distributions.

Quantifying uncertainty

After the distributions over the radiance field are learned, we are able to quantify the uncertainty of density or radiance of the scene. For each point x , we first sample S values from the distribution learned with normalizing flow and then use the sample variance to quantify the density uncertainty. The density uncertainty is a measure of convergence of the learned density. A similar procedure is employed to obtain

the radiance uncertainty. We can also quantify the uncertainty associated with the rendering image and depth map. Given any camera pose, we run the generator M times to sample M trajectories of random variables $(x^i, c^i)_{1:M}$ along each ray. Subsequently, we synthesize the images $\hat{I}_{1:M}$ and depth maps $\hat{D}_{1:M}$. The uncertainties can then be obtained through the color and depth variance of each pixel. Object edges often have high-density uncertainties due to the discontinuity of the depth whereas areas with high-frequency colors often have high radiance uncertainties. This kind of uncertainty is shown to be beneficial for active NeRF reconstruction in the experiments.

Training protocol

In contrast to standard NeRF, which primarily focuses on minimizing rendering loss, our GAN based on probabilistic NeRF incorporates a generator and discriminator that are trained in a min-max game framework. The generator comprises a deterministic branch and a probabilistic branch implemented with a normalizing flow. To ensure training stability, we adopt instance normalization (Ulyanov, Vedaldi, and Lempitsky 2016) and spectral normalization (Miyato et al. 2018) in the discriminator. To enhance convergence, we employ a small number of $M = 8$ samples for each point in the normalizing flow, gradually increasing it to 16 to generate M fake image patches. Simultaneously, we replicate the real patch M times before feeding it into the discriminator to maintain gradient balance during backpropagation. More implementation details can be found in the supplemental material.

Experiments

Datasets. We conducted evaluations of our method on three diverse publicly available datasets, namely LLFF (Mildenhall et al. 2021), Replica (Straub et al. 2019), and ScanNet (Dai et al. 2017). The LLFF dataset contains complex geometry and appearance, with camera views predominantly directed toward a single object. Notably, the LLFF dataset has also been utilized in S-NeRF (Shen et al. 2021) and CF-NeR (Shen et al. 2022). Additionally, we employed the Scannet and Replica datasets to evaluate the performance of our methods in challenging scenarios involving large scenes with diverse camera views. The Replica dataset comprises high-precision reconstructed scenes with nearly perfect rendering results, while the Scannet dataset consists of real-captured RGB-D sequences obtained using consumer-level cameras. Following the same experimental setting of (Yu et al. 2022), we utilized a subset of scenes in accordance with the methodology.

Metrics. To evaluate the predicted uncertainty and rendered quality, we employ standard metrics commonly used in the literature (Mildenhall et al. 2021; Shen et al. 2021, 2022). For uncertainty evaluation, we follow previous works by utilizing the AUSE (Area Under the Sparsification Error curve) metric proposed in (Bae, Budvytis, and Cipolla 2021; Poggi et al. 2020). AUSE evaluates the correlation between uncertainty and predicted error. Given an error metric (RMSE or MAE), we can sort the pixels of the uncertainty

Table 1: Comparison of uncertainty estimation and rendered color on the LLFF dataset. The best results are indicated in bold, with the second-best results underlined.

	Quality Metrics			Uncertainty Metrics	
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	AUSE RMSE \downarrow	AUSE MAE \downarrow
D.E. (Lakshminarayanan, Pritzel, and Blundell 2017)	22.32	0.788	0.236	0.0254	0.0122
Drop. (Gal and Ghahramani 2016)	21.90	0.758	0.248	0.0316	0.0162
NeRF-W (Martin-Brualla et al. 2021)	20.19	0.706	0.291	0.0268	0.0113
S-NeRF (Shen et al. 2021)	20.27	0.738	0.229	0.0248	0.0101
CF-NeRF (Shen et al. 2022)	21.96	<u>0.790</u>	<u>0.201</u>	<u>0.0177</u>	<u>0.0078</u>
FG-NeRF (Ours)	24.41	0.816	0.170	0.0118	0.0028

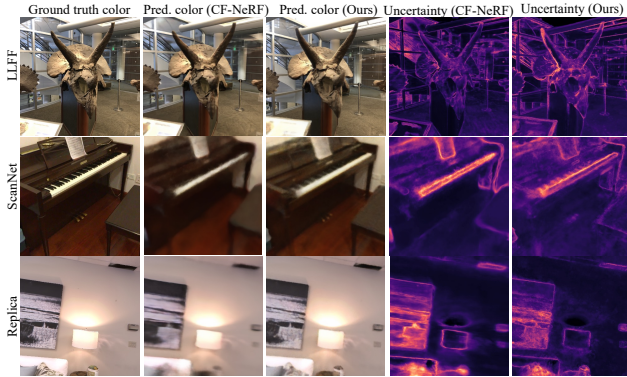


Figure 4: Quality comparison on the LLFF, ScanNet, and Replica datasets.

map and the error map respectively. Repeatedly removing the top 1% of the sorted error pixels, leads to a sparsification curve and an oracle curve. The AUSE is then computed as the area between these two curves. For rendered quality evaluation, we employ widely adopted metrics such as PSNR, SSIM (Structural Similarity Index), and LPIP (Learned Perceptual Image Patch Similarity). Detailed descriptions of these metrics can be found in the supplementary material.

Comparison on LLFF dataset

Following the same experimental setup in (Shen et al. 2022), we evaluate our uncertainty estimation method on the LLFF dataset to compare with existing methods (Lakshminarayanan, Pritzel, and Blundell 2017; Gal and Ghahramani 2016; Martin-Brualla et al. 2021; Shen et al. 2022, 2021). Note that, the camera views from LLFF dataset are primarily forward-facing to a central object within a small range of angles. Consequently, the coordinates can be constrained within the range of $[0, 1]$ using Normalized Device Coordinates (NDC), thereby facilitating the learning of the radiance field. The average results of eight scenes from LLFF dataset are presented in Tab. 1. Our method FG-NeRF achieves significant improvements in uncertainty and demonstrates better rendering quality as well (see Fig. 4 third rows for a visual comparison). Compared to the previous SOTA method CF-NeRF, the AUSE of our method decreased by 33.3% and 64.1% on RMSE and MAE, respectively, clearly demonstrating the superiority of our methods. Additionally, the results prove that the CF-NeRF outperforms the S-NeRF by relaxing the strong independence assumption. Furthermore, our method achieves the best performance by introducing

adversarial learning (Grover, Dhar, and Ermon 2018a) to sidestep any independence assumptions.

Comparison on Replica and Scannet dataset

Compared to object-level scenes, room-scale scenes present additional challenges due to their diverse camera views and complex geometry and appearance. In this context, we compare our methods with the previous state-of-the-art approach CF-NeRF (Shen et al. 2022) and NeurAR (Ran et al. 2022), specifically designed for scene-level uncertainty estimation. The results can be found in Table 2, and visual comparisons are illustrated in Figure 4. Notably, our method consistently achieves the best performance on both datasets. It is worth mentioning that CF-NeRF exhibits a performance drop when transitioning from the LLFF dataset (as shown in Table 1) to the ScanNet and Replica datasets (as shown in Table 2). In contrast, our method maintains a substantial level of performance, as observed in Table 1 on the LLFF dataset.

Additionally, we conduct a more challenging experiment to evaluate uncertainty estimation using a reduced number of frames (10% of the frames). The visualizations can be found in Figure 5, where it is evident that our method successfully predicts a clear boundary between the unobserved region and the observed space. Moreover, our approach consistently maintains accurate uncertainty estimation along the object edges.

Ablation study

We conduct an ablation study to verify the effectiveness of key techniques and hyperparameter settings. As shown in Tab. 3, the baseline is our full implementation with adversarial learning and has a deterministic branch. We first compare our proposed adversarial learning with plain L2 loss supervision as illustrated in Fig. 3. The L2 loss is the squared difference of the pixel color between rendered image patch \hat{P} and ground truth image patch P . In the experiments, we observed that both the rendering quality and the uncertainty metrics of adversarial learning are better than the case with plain L2 loss supervision. We also ablation studies the influence of the deterministic branch as proposed in the main pipeline. Experiments have shown a slight performance drop when the deterministic branch is removed. Additionally, we studied the different settings of the hyperparameters, e.g., the batch and image patch size, the sample size along each ray, and the scale annealing effect.

Table 2: Quantitatively comparison with baseline method on Scanet and Replica datasets.

Scene Type	Method	Quality Metrics			Uncertainty Metrics	
		PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	AUSE RMSE \downarrow	AUSE MAE \downarrow
ScanNet	CF-NeRF (Shen et al. 2022)	20.67	0.627	0.477	0.0293	0.0353
	NeurAR (Ran et al. 2023)	23.78	0.765	0.347	0.0265	0.0150
	Ours	24.76	0.793	0.388	0.0136	0.0043
Replica	CF-NeRF (Shen et al. 2022)	20.35	0.621	0.418	0.0196	0.0058
	NeurAR (Ran et al. 2023)	23.31	0.715	0.395	0.0175	0.0050
	Ours	25.40	0.816	0.364	0.0136	0.0040

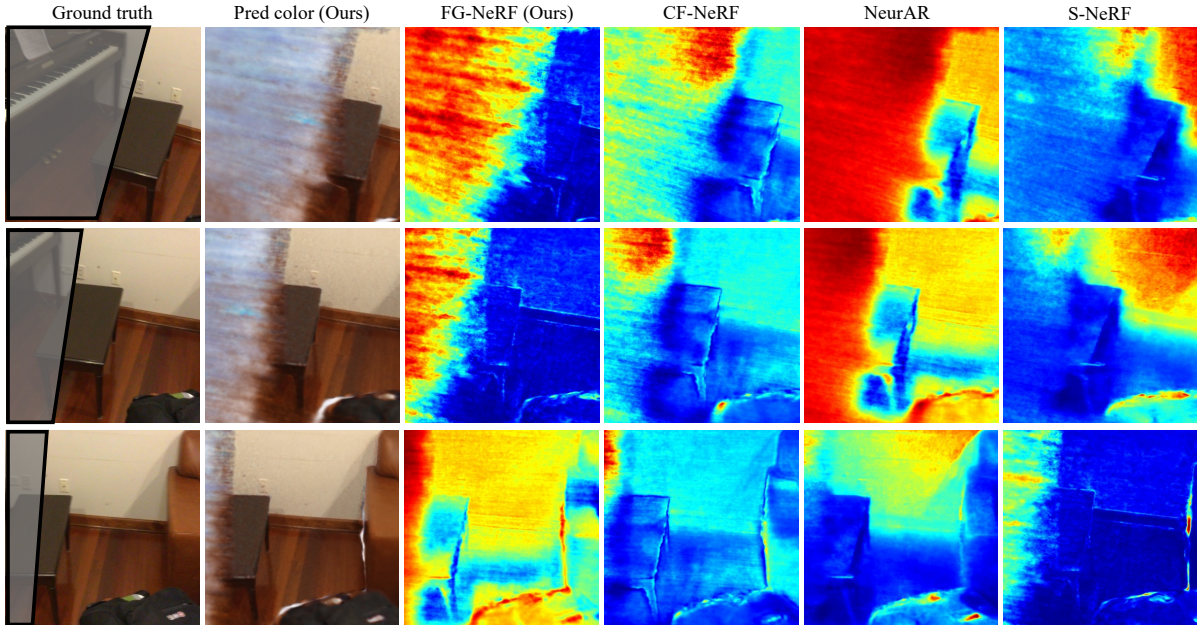


Figure 5: Comparison on uncertainty estimation under limited observations on ScanNet (the closer to blue, the smaller uncertainty). The unseen region is indicated by a gray mask.

Table 3: Ablations on LLFF Trex

		Quality Metrics			Uncertainty Metrics	
		PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	AUSE RMSE \downarrow	AUSE MAE \downarrow
Baseline		26.58	0.891	0.105	0.0070	0.0017
Color Supervision	L2 loss	25.89	0.842	0.197	0.0093	0.0043
Deterministic branch	off	24.96	0.723	0.285	0.0085	0.0049
Batch $B \times$ Patch K^2	4×16^2	25.43	0.817	0.091	0.0083	0.0024
Sample size M	24	25.44	0.854	0.146	0.0079	0.0019
	16	24.64	0.825	0.184	0.0082	0.0019
	8	23.58	0.777	0.248	0.0089	0.0022
Scale annealing	Off	24.38	0.814	0.199	0.0087	0.0021

Limitations

We acknowledge the presence of two main limitations in our approach. *Firstly*, there is a significant computational cost requirement. FG-NeRF, similar to S-NeRF and CF-NeRF, is trained to predict a radiance/density distribution, which is more challenging to learn compared to a simplistic radiance/density value. Consequently, these methods require more extensive computational resources than the original NeRF framework. Despite adopting a multi-level hashing coding method (Müller et al. 2022) for acceleration, our approach still takes nearly 6 GPU hours (NVIDIA V100) to obtain fully trained results on the LLFF dataset. This issue

can potentially be alleviated by integrating advanced NeRF acceleration techniques (Chen et al. 2022a,c). *Secondly*, the rendering quality of our method is relatively low when compared to the latest works in NeRF literature. We are keen on further enhancing the quality by leveraging scene prior information (Peng et al. 2020), incorporating advanced training strategies (Chen et al. 2022b), and employing novel backbone architectures (Sun et al. 2022).

Conclusions

In this paper, we propose FG-NeRF, a novel probabilistic neural radiance field, for high-quality uncertainty estimation. Our approach utilizes adversarial learning to model radiance/density distributions. Specifically, FG-NeRF performs as a generator and is trained in conjunction with a discriminator, forming an adversarial loss. By jointly supervising the colors of all pixels within an image patch, our method avoids making the independence assumption. Consequently, FG-NeRF effectively learns intricate radiance and density distributions and demonstrates its capacity to predict high-quality uncertainty, radiance, and density. Experi-

mental evaluations on both synthetic and real-world datasets showcase that our method outperforms existing methods, exhibiting lower rendering errors and providing more reliable uncertainty predictions. Moreover, we build and evaluate an active NeRF reconstruction framework based on FG-NeRF, showcasing the potential applications of this method.

References

- Bae, G.; Budvytis, I.; and Cipolla, R. 2021. Estimating and exploiting the aleatoric uncertainty in surface normal estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 13137–13146.
- Barron, J. T.; Mildenhall, B.; Tancik, M.; Hedman, P.; Martin-Brualla, R.; and Srinivasan, P. P. 2021. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 5855–5864.
- Blei, D. M.; Kucukelbir, A.; and McAuliffe, J. D. 2017. Variational inference: A review for statisticians. *Journal of the American statistical Association*, 112(518): 859–877.
- Brachmann, E.; Michel, F.; Krull, A.; Yang, M. Y.; Gumhold, S.; et al. 2016. Uncertainty-driven 6d pose estimation of objects and scenes from a single rgb image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 3364–3372.
- Cai, K.; Chen, W.; Wang, C.; Zhang, H.; and Meng, M. Q. 2021. Curiosity-based Robot Navigation under Uncertainty in Crowded Environments. *IEEE Robotics and Automation Letters*, 8: 800–807.
- Chan, E. R.; Lin, C. Z.; Chan, M. A.; Nagano, K.; Pan, B.; De Mello, S.; Gallo, O.; Guibas, L. J.; Tremblay, J.; Khamis, S.; et al. 2022. Efficient geometry-aware 3D generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 16123–16133.
- Chan, E. R.; Monteiro, M.; Kellnhofer, P.; Wu, J.; and Wetzstein, G. 2021. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 5799–5809.
- Chen, A.; Xu, Z.; Geiger, A.; Yu, J.; and Su, H. 2022a. Tensorf: Tensorial radiance fields. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXII*, 333–350. Springer.
- Chen, T.; Wang, P.; Fan, Z.; and Wang, Z. 2022b. Aug-NeRF: Training Stronger Neural Radiance Fields with Triple-Level Physically-Grounded Augmentations. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 15170–15181.
- Chen, Z.; Funkhouser, T. A.; Hedman, P.; and Tagliasacchi, A. 2022c. MobileNeRF: Exploiting the Polygon Rasterization Pipeline for Efficient Neural Field Rendering on Mobile Architectures. *ArXiv*, abs/2208.00277.
- Creswell, A.; White, T.; Dumoulin, V.; Arulkumaran, K.; Sengupta, B.; and Bharath, A. A. 2018. Generative adversarial networks: An overview. *IEEE signal processing magazine*, 35(1): 53–65.
- Dai, A.; Chang, A. X.; Savva, M.; Halber, M.; Funkhouser, T.; and Nießner, M. 2017. ScanNet: Richly-annotated 3D Reconstructions of Indoor Scenes. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE.
- Deng, Y.; Yang, J.; Xiang, J.; and Tong, X. 2022. Gram: Generative radiance manifolds for 3d-aware image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10673–10683.
- Gal, Y.; and Ghahramani, Z. 2016. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, 1050–1059. PMLR.
- Gao, K.; Gao, Y.; He, H.; Lu, D.; Xu, L.; and Li, J. 2022. NeRF: Neural Radiance Field in 3D Vision, A Comprehensive Review. *ArXiv*, abs/2210.00379.
- Geifman, Y.; Uziel, G.; and El-Yaniv, R. 2018. Bias-reduced uncertainty estimation for deep neural classifiers. *arXiv preprint arXiv:1805.08206*.
- Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2020. Generative adversarial networks. *Communications of the ACM*, 63(11): 139–144.
- Grover, A.; Dhar, M.; and Ermon, S. 2018a. Flow-gan: Combining maximum likelihood and adversarial learning in generative models. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.
- Grover, A.; Dhar, M.; and Ermon, S. 2018b. Flow-gan: Combining maximum likelihood and adversarial learning in generative models. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.
- Hendrycks, D.; Mazeika, M.; Kadavath, S.; and Song, D. 2019. Using self-supervised learning can improve model robustness and uncertainty. *Advances in neural information processing systems*, 32.
- Isola, P.; Zhu, J.-Y.; Zhou, T.; and Efros, A. A. 2017. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 1125–1134.
- Kirsch, W. 2019. An elementary proof of de Finetti’s theorem. *Statistics & Probability Letters*, 151: 84–88.
- Kosiorek, A. R.; Strathmann, H.; Zoran, D.; Moreno, P.; Schneider, R.; Mokrá, S.; and Rezende, D. J. 2021. Nerf-vae: A geometry aware 3d scene generative model. In *International Conference on Machine Learning*, 5742–5752. PMLR.
- Lakshminarayanan, B.; Pritzel, A.; and Blundell, C. 2017. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems*, 30.
- Lee, S.; Chen, L.; Wang, J.; Liniger, A.; Kumar, S.; and Yu, F. 2022. Uncertainty Guided Policy for Active Robotic 3D Reconstruction Using Neural Radiance Fields. *IEEE Robotics and Automation Letters*, 7(4): 12070–12077.
- Li, C.; Li, S.; Zhao, Y.; Zhu, W.; and Lin, Y. 2022a. RT-NeRF: Real-Time On-Device Neural Radiance Fields Towards Immersive AR/VR Rendering. *2022 IEEE/ACM In-*

- ternational Conference On Computer Aided Design (ICCAD), 1–9.
- Li, X.; Qiao, Y.-L.; Chen, P. Y.; Jatavallabhula, K. M.; Lin, M.; Jiang, C.; and Gan, C. 2023. PAC-NeRF: Physics Augmented Continuum Neural Radiance Fields for Geometry-Agnostic System Identification. *ArXiv*, abs/2303.05512.
- Li, Y.; Li, S.; Sitzmann, V.; Agrawal, P.; and Torralba, A. 2021. 3D Neural Scene Representations for Visuomotor Control. *ArXiv*, abs/2107.04004.
- Li, Z.; Li, L.; Ma, Z.; Zhang, P.; Chen, J.; and Zhu, J.-Z. 2022b. READ: Large-Scale Neural Scene Rendering for Autonomous Driving. *ArXiv*, abs/2205.05509.
- Malinin, A.; and Gales, M. 2018. Predictive uncertainty estimation via prior networks. *Advances in neural information processing systems*, 31.
- Martin-Brualla, R.; Radwan, N.; Sajjadi, M. S.; Barron, J. T.; Dosovitskiy, A.; and Duckworth, D. 2021. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 7210–7219.
- Meng, Q.; Chen, A.; Luo, H.; Wu, M.; Su, H.; Xu, L.; He, X.; and Yu, J. 2021. Gnerf: Gan-based neural radiance field without posed camera. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 6351–6361.
- Mildenhall, B.; Srinivasan, P. P.; Tancik, M.; Barron, J. T.; Ramamoorthi, R.; and Ng, R. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1): 99–106.
- Miyato, T.; Kataoka, T.; Koyama, M.; and Yoshida, Y. 2018. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*.
- Müller, T.; Evans, A.; Schied, C.; and Keller, A. 2022. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4): 1–15.
- Neff, T.; Stadlbauer, P.; Parger, M.; Kurz, A.; Mueller, J. H.; Chaitanya, C. R. A.; Kaplanyan, A.; and Steinberger, M. 2021. DONeRF: Towards Real-Time Rendering of Compact Neural Radiance Fields using Depth Oracle Networks. In *Computer Graphics Forum*, volume 40, 45–59. Wiley Online Library.
- Pan, X.; Lai, Z.; Song, S.; and Huang, G. 2022. ActiveNeRF: Learning Where to See with Uncertainty Estimation. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXIII*, 230–246. Springer.
- Peng, S.; Niemeyer, M.; Mescheder, L. M.; Pollefeys, M.; and Geiger, A. 2020. Convolutional Occupancy Networks. *ArXiv*, abs/2003.04618.
- Poggi, M.; Aleotti, F.; Tosi, F.; and Mattoccia, S. 2020. On the uncertainty of self-supervised monocular depth estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 3227–3237.
- Ran, Y.; Zeng, J.; He, S.; Chen, J.; Li, L.; Chen, Y.; Lee, G.; and Ye, Q. 2023. NeurAR: Neural Uncertainty for Autonomous 3D Reconstruction with Implicit Neural Representations. *IEEE Robotics and Automation Letters*.
- Ran, Y.; Zeng, J.; He, S.; Li, L.; Chen, Y.; Lee, G.; Chen, J.; and Ye, Q. 2022. NeurAR: Neural Uncertainty for Autonomous 3D Reconstruction. *arXiv preprint arXiv:2207.10985*.
- Schwarz, K.; Liao, Y.; Niemeyer, M.; and Geiger, A. 2020. Graf: Generative radiance fields for 3d-aware image synthesis. *Advances in Neural Information Processing Systems*, 33: 20154–20166.
- Shen, J.; Agudo, A.; Moreno-Noguer, F.; and Ruiz, A. 2022. Conditional-flow nerf: Accurate 3d modelling with reliable uncertainty quantification. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part III*, 540–557. Springer.
- Shen, J.; Ruiz, A.; Agudo, A.; and Moreno-Noguer, F. 2021. Stochastic neural radiance fields: Quantifying uncertainty in implicit 3d representations. In *2021 International Conference on 3D Vision (3DV)*, 972–981. IEEE.
- Straub, J.; Whelan, T.; Ma, L.; Chen, Y.; Wijmans, E.; Green, S.; Engel, J. J.; Mur-Artal, R.; Ren, C.; Verma, S.; Clarkson, A.; Yan, M.; Budge, B.; Yan, Y.; Pan, X.; Yon, J.; Zou, Y.; Leon, K.; Carter, N.; Briales, J.; Gillingham, T.; Mueggler, E.; Pesqueira, L.; Savva, M.; Batra, D.; Strasdat, H. M.; Nardi, R. D.; Goesele, M.; Lovegrove, S.; and Newcombe, R. 2019. The Replica Dataset: A Digital Replica of Indoor Spaces. *arXiv preprint arXiv:1906.05797*.
- Sun, J.; Chen, X.; Wang, Q.; Li, Z.; Averbuch-Elor, H.; Zhou, X.; and Snavely, N. 2022. Neural 3D Reconstruction in the Wild. *ACM SIGGRAPH 2022 Conference Proceedings*.
- Teye, M.; Azizpour, H.; and Smith, K. 2018. Bayesian uncertainty estimation for batch normalized deep networks. In *International Conference on Machine Learning*, 4907–4916. PMLR.
- Ulyanov, D.; Vedaldi, A.; and Lempitsky, V. 2016. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*.
- Vora, S.; Radwan, N.; Greff, K.; Meyer, H.; Genova, K.; Sajjadi, M. S. M.; Pot, E.; Tagliasacchi, A.; and Duckworth, D. 2021. NeSF: Neural Semantic Fields for Generalizable Semantic Segmentation of 3D Scenes.
- Wang, P.; Liu, L.; Liu, Y.; Theobalt, C.; Komura, T.; and Wang, W. 2021. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689*.
- Winkler, C.; Worrall, D.; Hooeboom, E.; and Welling, M. 2019. Learning likelihoods with conditional normalizing flows. *arXiv preprint arXiv:1912.00042*.
- Yang, G.; Huang, X.; Hao, Z.; Liu, M.-Y.; Belongie, S.; and Hariharan, B. 2019. Pointflow: 3d point cloud generation with continuous normalizing flows. In *Proceedings of the IEEE/CVF international conference on computer vision*, 4541–4550.
- Yariv, L.; Gu, J.; Kasten, Y.; and Lipman, Y. 2021. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems*, 34: 4805–4815.

Yu, Z.; Peng, S.; Niemeyer, M.; Sattler, T.; and Geiger, A. 2022. MonoSDF: Exploring Monocular Geometric Cues for Neural Implicit Surface Reconstruction. *ArXiv*, abs/2206.00665.

Zhang, K.; Riegler, G.; Snavely, N.; and Koltun, V. 2020. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*.

Zhi, S.; Laidlow, T.; Leutenegger, S.; and Davison, A. J. 2021. In-place scene labelling and understanding with implicit scene representation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 15838–15847.