

Aug-NeRF: Training Stronger Neural Radiance Fields with Triple-Level Physically-Grounded Augmentations

Tianlong Chen^{1*}, Peihao Wang^{1*}, Zhiwen Fan¹, Zhangyang Wang¹
¹The University of Texas at Austin

{tianlong.chen, peihaowang, zhiwenfan, atlaswang}@utexas.edu

Abstract

Neural Radiance Field (NeRF) regresses a neural parameterized scene by differentially rendering multi-view images with ground-truth supervision. However, when interpolating novel views, NeRF often yields inconsistent and visually non-smooth geometric results, which we consider as a **generalization gap** between seen and unseen views. Recent advances in convolutional neural networks have demonstrated the promise of advanced robust data augmentations, either random or learned, in enhancing both in-distribution and out-of-distribution generalization. Inspired by that, we propose Augmented NeRF (**Aug-NeRF**), which for the first time brings the power of robust data augmentations into regularizing the NeRF training. Particularly, our proposal learns to seamlessly blend worst-case perturbations into three distinct levels of the NeRF pipeline with physical grounds, including (1) the input coordinates, to simulate imprecise camera parameters at image capture; (2) intermediate features, to smoothen the intrinsic feature manifold; and (3) pre-rendering output, to account for the potential degradation factors in the multi-view image supervision. Extensive results demonstrate that Aug-NeRF effectively boosts NeRF performance in both novel view synthesis (up to 1.5dB PSNR gain) and underlying geometry reconstruction. Furthermore, thanks to the implicit smooth prior injected by the triple-level augmentations, Aug-NeRF can even recover scenes from heavily corrupted images, a highly challenging setting untackled before. Our codes are available in <https://github.com/VITA-Group/Aug-NeRF>.

1. Introduction

Neural radiance fields (NeRF) [29] and its variants have demonstrated impressive progresses in learning to represent 3D objects and scenes from images towards photo-realistic novel view synthesis. NeRF leverages a multi-layer perceptron (MLP) to implicitly modeling the mapping from an input 5D coordinates (i.e., 3D coordinates (x, y, z) and 2D viewing

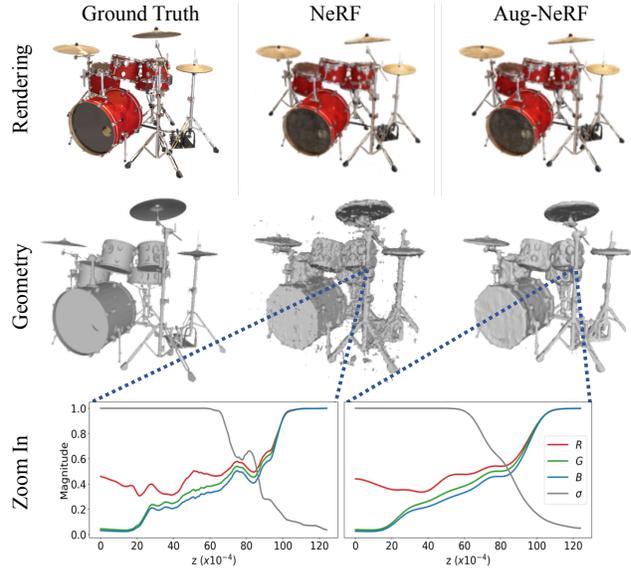


Figure 1. Comparisons between Aug-NeRF (ours) and NeRF [29]. From upper to bottom, we present the test-set synthesized views, 3D geometry, and their zoom-in RGB σ distributions, respectively.

directions (θ, ϕ)) to volume density σ and view-dependent emitted radiance color (r, g, b) at the corresponding position in the scene. Then, the obtained continuous 5D function (i.e., MLP) can be utilized to generate novel views with traditional volume rendering mechanisms.

Although NeRF is capable of producing novel views, it unfortunately suffers from inconsistent and non-smooth geometries since the vanilla MLP lacks geometry-awareness. For example, as shown in Fig. 1, the depth maps and 3D geometries of the scene generated by NeRF show obvious discontinuity and outliers, especially around the edge of objects. Considering that the quality of reconstructed geometry plays a central role in view rendering, that might account for NeRF’s limited generalization to unseen views.

To fill in this research gap, a straightforward solution is introducing explicit geometric regularizers like Laplacian [17, 47] or total variation (TV) [88] to enhance the continuity. However, these explicit regularizers are often

*Equal Contribution.

found to constrain the representation flexibility of MLP too aggressively, resulting in inferior performance. Recent advances in robust data augmentations [75] establish promising successes in image recognition in terms of both improved functional smoothness and generalization.

Motivated by that, we design an **Augmented NeRF (Aug-NeRF)** training framework, which injects worst-case perturbations [23] to implicitly regularize the NeRF pipeline with physical foundations. Specifically, Aug-NeRF considers to regularize three different levels, including (i) *the input coordinates*, where perturbations can imitate the inaccurate camera poses during collecting images; (ii) *the intermediate features*, in order for a smooth/flat model loss landscape [5, 8, 73] when fitting objects’ 3D geometries that is believed to enhance generalization; (iii) *the pre-rendering output*, to model potential degradation factors in the image supervision. As presented in Fig. 1, our Aug-NeRF achieves smoother and more consistency reconstructed geometry and improved unseen view synthesis. Additionally, we find Aug-NeRF to show surprising resilience towards severely corrupted supervision images. The main contributions of this paper can be summarized as follows:

- We reveal the existence of highly non-smooth geometries in representing scenes as neural radiance fields (NeRF), which we regard as a crucial bottleneck of NeRF’s generalization ability to unseen views.
- To address such limitation of NeRF, we propose Aug-NeRF, a triple-level, physically-grounded augmented training pipeline, by leveraging worst-case perturbations to implicated regularize the input coordinate, intermediate feature, and pre-rendering output levels.
- Extensive experiments validate the effectiveness of our proposal on diverse scene synthesis tasks, to endow NeRF with smoothness-aware geometry reconstruction, enhanced generalization to synthesizing unseen views, and stronger tolerance of noisy supervisions.

2. Related Work

Adversarial Training and Robust Augmentation. It is well-known that deep networks are vulnerable to imperceptible worst-case perturbations [11, 16, 23]. Numerous defense mechanisms [31, 40, 49, 57, 59, 84] have been invented to address the issue, where adversarial training (AT) approaches [11, 16, 23] remains as the *de-facto*. Although conventional AT enhances model robustness at the price of compromising the standard accuracy [64], recent studies reveal AT can be harnessed to enhance models’ standard generalization as well [9, 65, 71, 75, 90]. Taking [75] for example, it applies adversarial perturbations to input samples as a form of data augmentation, and shows to improve image classification on the clean dataset. [9, 65, 90] apply worst-case perturbations to the input embedding for natural language

understanding, language modeling, and vision-and-language tasks, all successfully boosting their standard generalization. [14, 42, 66, 87] constructed more sophisticated variations of robust augmentations, including both data-driven and heuristic components, to improve model generalization further. However, such robust augmentations on inputs or intermediate features, to our best knowledge, have not been studied in the view synthesis field. This paper explores this possibility by looking into the intrinsic physical grounds.

Neural 3D Representations. Classic 3D reconstruction approaches utilizes discrete representations such as point clouds [1, 74], meshes [45, 46, 63], multi-plane images [28, 55, 56, 89], depth maps [12, 51, 77] and voxel grids [18, 21, 39, 50, 53, 60]. Neural implicit representations leverage coordinate-based neural networks to approximate visual signals [27, 35, 38]. Such ideas have been successfully applied to both 2D images [20, 52, 62] and 3D objects [6, 48, 52]. Recent advances follow differentiable rendering and end-to-end optimization to reconstruct the neural 3D scene from 2D image supervision [29, 33, 79]. Liu *et al.* [20] presented the first usage of neural implicit function to infer 3D representation with differentiable rendering. DVR [33] and IDR [79] adopt surface rendering to reconstruct implicit iso-surface by supervising on both images and pixel-accurate object masks.

NeRF [29] pioneered to use differentiable volumetric rendering to optimize a neural radiance field, and achieved more photorealistic and view-consistent results. Many works continue to improve its training and rendering accuracy, efficiency, and generalization. NeRF++ [85] separates two NeRFs to handle foreground and background, respectively. NeRF-W [24] tackles unstructured photos via modeling transient noises and uncertainty. MipNeRF [3] mitigates objectionable aliasing artifacts for NeRF to represent fine details. HyperNeRF [36] introduces topology-aware level-set methods to rectify NeRF geometry especially for dynamics. [10, 41, 44, 72] extend NeRF with lighting and rendering modeling. [34, 67, 78] enhance the underlying geometries reconstructed by NeRF by adopting surface representation in the place of the density volume. [4, 68, 83] leverage multi-view spatial image feature or semi-reconstructed 3D information to reduce input view number and enable generalization to new scenes. [26, 70, 80] free NeRF from accurate camera pose estimation. Acceleration of NeRF training and inference have also been discussed in [43, 61, 82]. Despite so many exciting progresses, studying NeRF’s training stability and data robustness remains an open question.

3. Preliminaries

NeRF models the underlying 3D scene as a continuous volumetric radiance field of color and density. Formally, a typical radiance field can be written as $F : (\mathbf{x}, \boldsymbol{\theta}) \mapsto (\mathbf{c}, \sigma)$, where $\mathbf{x} \in \mathbb{R}^3$ is the spatial coordinate, $\boldsymbol{\theta} \in [-\pi, \pi]^2$ indicates the view direction, and $\mathbf{c} \in \mathbb{R}^3, \sigma \in \mathbb{R}_+$ represent

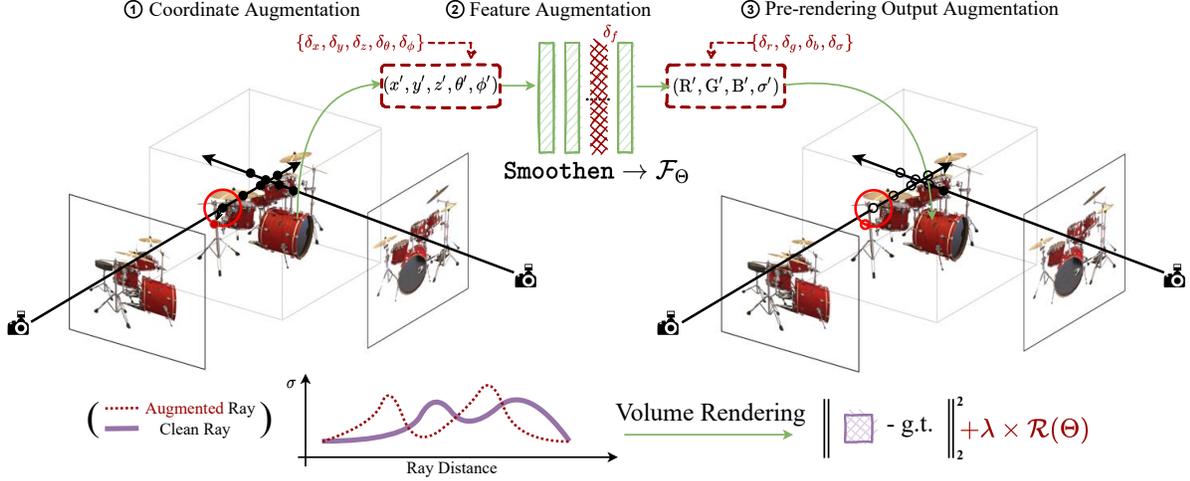


Figure 2. The overall pipeline of our proposed Aug-NeRF. The worst-case perturbations are generated at the triple levels of the NeRF pipeline: ① input coordinates, ② intermediate features, and ③ pre-rendering output.

the RGB color and density, respectively. NeRF further parameterizes this 5D-valued function by a composition of Positional Embedding (PE) and the MLP $F_{\Theta} = \gamma \circ \text{MLP}_{\Theta}$, where γ is a Fourier feature mapping network [62], Θ is the network weights. Given a radiance field, NeRF follows the classical volume rendering to render an arbitrary view [25].

Our goal is to fit a neural radiance from calibrated RGB images captured from multiple views. Suppose we have a set of images with corresponding extrinsic parameters. NeRF simulates the physical imaging process, by casting a ray $\mathbf{r} = (\mathbf{o}, \mathbf{d}, \boldsymbol{\theta})$ for each pixel via inverse perspective projection with respect to the camera pose, where $\mathbf{o} \in \mathbb{R}^3$ denotes the optical center of camera, $\mathbf{d} \in \mathbb{R}^3$ is the direction of the ray, and $\boldsymbol{\theta} \in [-\pi, \pi]^2$ is the angular view direction (see Fig. 2). We collect all pairs of rays and pixel colors as the training set $\mathcal{R} = \{(\mathbf{r}_i, \hat{\mathbf{C}}_i)\}_{i=1}^N$, where N is the total number of rays, and $\hat{\mathbf{C}}_i$ denotes the ground-truth color of the i -th ray. To simulate the color of a ray, NeRF first partitions K evenly-spaced bins between the near-far bound $[t_n, t_f]$ along the ray, and then uniformly samples one point within each bin: $t_k \sim \mathcal{U}[t_n + (k-1)(t_f - t_n)/K, t_f + k(t_f - t_n)/K]$. Afterwards, NeRF numerically evaluates volumetric ray integration [25] via the following equation:

$$\mathbf{C}(\mathbf{r}|\Theta) = \sum_{k=1}^K T(k)(1 - \exp(-\sigma_k \Delta t_k)) \mathbf{c}_k$$

$$\text{where } T(k) = \exp\left(-\sum_{l=1}^{k-1} \sigma_l \Delta t_l\right), \quad (1)$$

where $\Delta t_k = t_{k+1} - t_k$, and $(\mathbf{c}_k, \sigma_k) = F_{\Theta}(\mathbf{o} + t_k \mathbf{d}, \boldsymbol{\theta})$. With this forward model, NeRF optimizes the expected L_2 distance between rendered ray colors and ground-truth pixel colors as follows:

$$\mathcal{L}(\Theta|\mathcal{R}) = \mathbb{E}_{(\mathbf{r}, \hat{\mathbf{C}}) \sim \mathbb{P}(\mathcal{R})} \left\| \mathbf{C}(\mathbf{r}|\Theta) - \hat{\mathbf{C}} \right\|_2^2, \quad (2)$$

where $\mathbb{P}(\cdot)$ defines a probability measure supported in the ray space \mathcal{R} .

4. Methodology

Overview. NeRF conducts uniform sampling along each ray and interpolates a continuous radiance field via an MLP. However, we argue that the point sampling and the MLP interpolation can never be optimal during training dynamics due to the biased sampling strategy and non-smoothness of MLP. To this end, we propose to train NeRF with a smoothing prior. Sec. 4.1 provides a probabilistic interpretation of this intuition. Different from explicit smoothness modeling, e.g., total variation penalty or low rank prior, we utilize worst-case perturbations as a data-adaptive regularization. We call this training strategy Aug-NeRF.

An overview of our Aug-NeRF is presented in Fig. 2. Following the rendering pipeline of NeRF, Aug-NeRF injects adversarial noises into the following stages: point sampling, intermediate features, and MLP outputs. Each perturbation is searched within a small range to maximize the final loss. It could be treated as a regularization to be jointly minimized with the original training loss (see Sec. 4.2).

4.1. NeRF as Maximum A Posterior

Fitting a neural radiance field to satisfy multi-view observations can be modeled as a Maximal Likelihood (ML) problem $\Theta^* = \arg \max_{\Theta} \mathbb{P}(\mathcal{R}|\Theta)$, which can be derived as:

$$\Theta^* = \arg \min_{\Theta} -\mathbb{E}_{(\mathbf{r}, \hat{\mathbf{C}}) \sim \mathbb{P}(\mathcal{R})} \log \mathbb{P}(\mathbf{r}, \hat{\mathbf{C}}|\Theta),$$

by assuming each ray is conditionally independent given network parameters. Optimizing NeRF by MSE loss (Eqn. 2) can be obtained by regarding the conditional distribution

$\mathbb{P}(\mathbf{r}, \widehat{\mathbf{C}}|\Theta)$ as a Gaussian distribution:

$$\mathbb{P}(\mathbf{r}, \widehat{\mathbf{C}}|\Theta) = \frac{1}{Z} \exp\left(-\frac{1}{2\Sigma} \left\| \mathbf{C}(\mathbf{r}|\Theta) - \widehat{\mathbf{C}} \right\|_2^2\right),$$

where Z is a normalization term, and Σ is the variance.

However, the maximum likelihood does not introduce any prior on the reconstructed NeRF as MLP is a universal approximator. Instead, we consider the Maximum A Posterior (MAP) form $\mathbb{P}(\Theta|\mathcal{R})$ to inject the prior for robust training. By Bayesian rule, we have $\mathbb{P}(\Theta|\mathcal{R}) \propto \mathbb{P}(\mathcal{R}|\Theta) \mathbb{P}(\Theta)$, where $\mathbb{P}(\Theta)$ is some prior distribution of the network weights Θ . Hence, maximizing this posterior probability is equivalent to minimizing the original loss (Eqn. 2) plus a penalty term:

$$\mathcal{L}(\Theta|\mathcal{R}) = \mathbb{E}_{(\mathbf{r}, \widehat{\mathbf{C}}) \sim \mathbb{P}(\mathcal{R})} \left\| \mathbf{C}(\mathbf{r}|\Theta) - \widehat{\mathbf{C}} \right\|_2^2 + \lambda R(\Theta), \quad (3)$$

where $R(\Theta) = -\log \mathbb{P}(\Theta)/\lambda$. Here we expect Θ to induce a geometry-aware smooth F_Θ .

4.2. Regularize NeRF with Robust Augmentations

Imposing smoothness onto NeRF can be done in many explicit ways, such as regularizing total variation [88], Laplacian of surface [7, 54, 81], etc. However, those regularizers are often not sufficiently data-adaptive, and can constrain the representation flexibility too aggressively, as evidenced in Sec. 5.3. Also, their computation also usually operates on discretized volumetric representations, and needs extra differentiation steps to be added in NeRF.

Recent works [5, 9, 65, 71, 75, 90] suggest a promising alternative by integrating worst-case adversarial perturbations as data augmentations (i.e., AT). AT restricts the change of loss when its input is perturbed, leading to flattening the loss landscape [30, 58]. As a result, the trained network’s intrinsic feature manifold and loss landscape become smoother. Prevailing theories [15, 19, 32] link the generalization ability of deep networks to the geometry of the loss landscape; in particular, a model trained to converge to wide valleys (i.e., flat basins) in loss landscape shows better generalization ability as well as robustness to distributional shifts.

NeRF is trained by given 2D image views (often with known camera poses) and is tested to synthesize novel views from unseen angles. Intuitively, the unsatisfactory novel view synthesis could be seen as a training-testing “generalization gap” issue. This inspires us to incorporate robust augmentations into NeRF to induce a data-adaptive smoothness prior that enhances generalization.

Designing dedicated perturbations for NeRF is far from trivial due to its inherent physics. Unlike conventional deep models, the forward pass of NeRF consists of two white-box simulating stages (point sampling, volumetric rendering) and one black-box network mapping stage. We propose to inject worst-case perturbations into all three levels: coordinates,

intermediate features of MLP, and pre-rendering MLP output: all with clear physical meanings. Formally, our approach can be formulated as a min-max game:

$$\min_{\Theta} \mathbb{E}_{(\mathbf{r}, \widehat{\mathbf{C}}) \sim \mathbb{P}(\mathcal{R})} \max_{\delta} \left\| \mathbf{C}^\dagger(\mathbf{r}|\Theta, \delta) - \widehat{\mathbf{C}} \right\|_2^2, \quad (4)$$

where $\delta = (\delta_p, \delta_f, \delta_r) \in \mathcal{S}_p \times \mathcal{S}_f \times \mathcal{S}_r$,

where δ_p , δ_f , and δ_o are the perturbations to be learned and injected to the input coordinate, intermediate MLP feature, and pre-rendering RGB- σ output, respectively, where $\mathcal{S}_p \subseteq \mathbb{R}^6$, $\mathcal{S}_f \subseteq \mathbb{R}^D$, and $\mathcal{S}_r \subseteq \mathbb{R}^4$ are the corresponding perturbation search range, D is the hidden dimension of the MLP. We elaborate on each perturbation as below.

Input Coordinate Perturbation. The original NeRF first randomly samples point along each ray and then conducts importance sampling to simulate the quadrature of the integration. This strategy also mitigates overfitting and produces smoother scene representation [29]. Arandjelovic *et al.* [2] further proposes an attention-guided sampling scheme to refine this process. However, our insight is that using either coarse-to-fine or learning-based sampling will cause the sampling to overfit the density distribution of the currently rendered ray, which might hold back NeRF when the density field is biased or cannot generalize.

To this end, we propose to produce a worst-case point sampling during training, to simulate a test-time “distributional shift” for NeRF to handle. To be specific, we search a coordinate perturbation $\delta_{xyz} = (\delta_x, \delta_y, \delta_z)$ following Eqn. 4. The coordinate perturbation $\delta_p = (\delta_t, \delta_{xyz}, \delta_\theta)^T$ consists of three parts: 1) the along-ray perturbation $\delta_t \in \mathbb{R}$ shifts point samples along the ray, 2) the point position perturbation $\delta_{xyz} \in \mathbb{R}^3$ is added to the direct input of the NeRF MLP, 3) in addition, we also inject the perturbation $\delta_\theta \in \mathbb{R}^3$ to the view direction. Formally, given the perturbation δ_p , the input of MLP turns out to be:

$$t_k^\dagger = t_k + \delta_t, \quad \theta^\dagger = \theta + \delta_\theta, \quad \mathbf{p}_k^\dagger = \mathbf{o} + t_k^\dagger \mathbf{d} + \delta_{xyz}.$$

The constraint set for δ_t is defined as $\delta_t \leq |\alpha_t(t_{k+1} - t_k)|$, where α_t is a hyperparameter. The coordinate perturbation δ_{xyz} lies in a ball $\mathcal{B}(0, \epsilon_p)$ to constrain points with a cylinder along the ray. View direction perturbation δ_θ is restricted within the conical frustum $[-\epsilon_p/2f, \epsilon_p/2f]^2$, where f is the focal length, and ϵ_p is the pixel size.

Pre-Rendering Output Perturbation. NeRF next maps points on a ray to the corresponding color and density, then conducts volumetric rendering to compose these point values into the 2D pixel values. As shown by Fig. 1, the reconstructed shape can be noisy and discontinuous. We attribute these artifacts to two reasons: (i) neural implicit functions represented by MLP are not necessarily smooth [62, 76]. When zooming in, we observe the function landscape to be



Figure 3. Comparisons on the test-set views for scenes from the realistic LLFF dataset [29]: local zoom-in in the red box.

rugged; (ii) the MLP output goes through volumetric rendering to form the RGB output. As the volumetric rendering itself has smoothing effects owing to its point-by-point accumulation, it might “mask” the non-smoothness and noise of the pre-rendering results hence they cannot be effectively eliminated at supervised training.

Inspired by robust training enhancing output smoothness [5, 9, 65, 71, 75, 90], we propose to intentionally corrupt the output of the MLP with worst-case perturbation, in order to encourage the output smoothness of the MLP, which in turn smooths the NeRF underlying geometry. Given the pre-rendering perturbation $\delta_r = (\delta_c, \delta_\sigma)$, $\delta_c = (\delta_r, \delta_g, \delta_b)$, we perturb the rendering in Eqn. 1 by:

$$C^\dagger(\mathbf{r}|\Theta) = \sum_{k=1}^K T(k)(1 - \exp(-(\sigma_k + \delta_\sigma)\Delta t_k^\dagger))(c_k + \delta_c),$$

where $(c_k, \sigma_k) = F_\Theta(\mathbf{p}_k^\dagger, \theta^\dagger)$ are outputs by perturbed coordinates, $T(k) = \exp\left(-\sum_{l=1}^{k-1}(\sigma_l + \delta_\sigma)\Delta t_l^\dagger\right)$ is the transmittance term, $\Delta t_k^\dagger = t_{k+1}^\dagger - t_k^\dagger$ is the interval of integral, and δ_c, δ_σ correspond to color and density perturbations, respectively. We fix the constraint set as $[-\epsilon_c, \epsilon_c]^3 \times [-\epsilon_\sigma, \epsilon_\sigma]$. δ_c, δ_σ will be further clamped to make sure c_k, σ_k lie between $[0, 1]$.

Intermediate Feature Perturbation. In addition to perturbing per-rendering color and ray density, we also inject adversarial noise into the intermediate features. As revealed by [5], augmenting intermediate features can further smooth learned functional mappings, more than just augmenting inputs or outputs. To be specific, according to [29], the backbone MLP can be written as $(c(\mathbf{p}), \sigma(\mathbf{p})) = F_\Theta(\mathbf{p}, \theta) =$

$(g \circ f(\mathbf{p}, \theta), h \circ f(\mathbf{p}))$, where $f(\cdot)$ (with positional encoding) maps a coordinate to a D -dimension feature vector, and $g(\cdot), h(\cdot)$ project it to RGB color and density, respectively. We crafted the worst-case perturbations as follows:

$$c_k = g(f(\mathbf{p}_k^\dagger) + \delta_f, \theta^\dagger), \quad \sigma_k = h(f(\mathbf{p}_k^\dagger) + \delta_f).$$

Intermediate feature perturbation is searched over $\mathcal{S}_f = [-\epsilon_f, \epsilon_f]^D$ with a hyperparameter ϵ_f . We also test various injection points of the backbone MLP in Sec. 5.3.

4.3. Optimization

To search for the worst-case perturbation in Eqn. 4, we introduce a theoretically guaranteed way to reach the maximum. We only consider additive perturbation here, and all search spaces (i.e., $\mathcal{S}_p, \mathcal{S}_f, \mathcal{S}_r$) are defined as ℓ_p norm ball with a radius $\epsilon > 0$. The radius ϵ is the maximum magnitude of the perturbation, which can roughly signify the strength of the perturbation. The perturbations can be accurately estimated by multi-step Projected Gradient Descent (PGD). Taking ℓ_∞ norm ball for example:

$$\delta^{(t+1)} = \Pi_{\|\delta\|_\infty \leq \epsilon} \left[\delta^{(t)} + \alpha \cdot \text{sgn}(\nabla_\delta \mathcal{L}(\Theta|\mathcal{R}, \delta)) \right] \quad (5)$$

where α is the step size of the inner maximization, $\Pi[\cdot]$ denotes a projection operator, $\text{sgn}(\cdot)$ takes the sign of the input, and $\mathcal{L}(\Theta|\mathcal{R}, \delta)$ represents the MSE loss between perturbed color C^\dagger and ground-truth color \hat{C} (see Eqn. 4).

After incorporating all augmentations, the full training objective is defined as ($\lambda = 0.1$ as tuned by grid search):

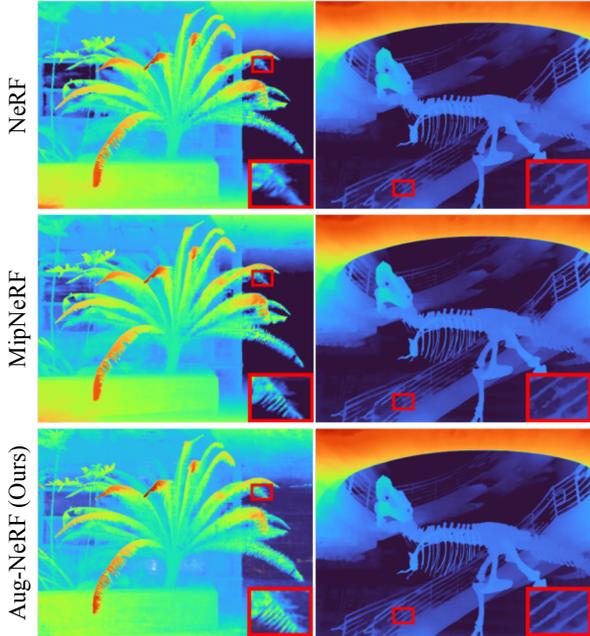


Figure 4. Comparisons of learned depth maps on scenes from LLFF dataset. The local zoom-in is placed in the red box.

$$\mathbb{E}_{(r, \hat{C}) \sim \mathbb{P}(\mathcal{R})} \left[\underbrace{\|C(r|\Theta) - \hat{C}\|_2^2}_{\text{photometric loss}} + \lambda \max_{\delta} \underbrace{\|C^\dagger(r|\Theta, \delta) - \hat{C}\|_2^2}_{\text{adversarial reg.}} \right].$$

5. Experiments

5.1. Implementation details.

Datasets. We evaluate our proposals on public representative datasets of both LLFF [28] and NeRF-Synthetic [29]. Particularly, the face-forwarding scenes {“fern”, “orchids”, “trex”} from LLFF dataset and {“drums”, “ship”, and “chair”} instances in 360° NeRF-Synthetic dataset are adopted in our experiments. To accelerate training, we down-sampled LLFF dataset by 1/8 and 360° NeRF-Synthetic dataset by 1/2.

Training. We employ the same MLP architecture and training recipe with the original NeRF. Aug-NeRF is trained for 500K iterations to guarantee convergence. All hyperparameters are carefully tuned by a grid search and the best configuration is applied to all experiments, as demonstrated in Sec. 5.3. NeRF models are trained on a NVIDIA RTX A6000 GPU with 48 GB memory.

Evaluation. We report three error metrics including peak signal-to-noise ratio (PSNR), the structural similarity index measure (SSIM) [69], and learned perceptual image patch similarity (LPIPS) [86]. Meanwhile, to provide a comprehensive comparison, we also follow [3] and show an “average” error metric by computing the geometric mean of

Table 1. Quantitative comparison of our Aug-NeRF against NeRF and other top-performers for novel view synthesis. Performance is reported on the LLFF test set. \uparrow/\downarrow means that larger/smaller numbers denote better performance.

Scene “fern”	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Average \downarrow
SRN [52]	21.37	0.611	0.459	0.128
LLFF [52]	22.85	0.753	0.247	0.086
NeRF [29]	25.17	0.792	0.280	0.073
MipNeRF [3]	26.24	0.839	0.193	0.057
Aug-NeRF (Ours)	26.51	0.830	0.168	0.054
Scene “orchids”	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Average \downarrow
SRN [52]	17.37	0.611	0.467	0.175
LLFF [28]	18.52	0.588	0.313	0.141
NeRF [29]	20.36	0.641	0.321	0.121
MipNeRF [3]	20.87	0.663	0.262	0.108
Aug-NeRF (Ours)	21.60	0.675	0.243	0.099
Scene “trex”	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Average \downarrow
SRN [52]	22.87	0.761	0.298	0.091
LLFF [28]	24.15	0.857	0.222	0.069
NeRF [29]	26.80	0.880	0.249	0.056
MipNeRF [3]	27.55	0.894	0.208	0.049
Aug-NeRF (Ours)	28.17	0.881	0.206	0.048

$\text{MSE} = 10^{-\frac{\text{PSNR}}{10}}$, $\sqrt{1 - \text{SSIM}}$, and LPIPS.

Baseline and Comparison Variants. Our Aug-NeRF is established on the vanilla NeRF [29]. Two groups of current top-performers for view synthesis are compared, including (i) NeRF-based approaches: NeRF [29] and MipNeRF [3]; and (ii) classical methods: Neural Volume (NV) [21], Scene Representation Network (SRN) [52], and Local Light Field Fusion (LLFF) [28]. For a fair comparison, all above models are trained/tested on the same views of identical scenes.

5.2. Improved NeRF with Augmentations

Results on LLFF and 360° NeRF-Synthetic datasets. In this section, we validate our proposed Aug-NeRF on LLFF and 360° NeRF-Synthetic datasets across six representative scenes. Quantitative comparisons against vanilla NeRF and other top-performing algorithms like {MipNeRF [3], NV [21], SRN [52], LLFF [28]} are provided in Tab. 1 and 2, together with qualitative test views presented in Fig. 3. These results convey several observations:

- ① Aug-NeRF reduces average error by 14.3% \sim 26.0% and 12.5% \sim 44.7% on the LLFF and 360° NeRF-Synthetic datasets, respectively. It consistently outperforms NeRF on all metrics by a large margin, e.g., {1.34, 1.24, 1.37, 1.33, 0.53, 0.87} PSNR improvements at scenes {“fern”, “orchids”, “trex”, “drums”, “ship”, “chair”}, showing impressive “generalization” boosts on unseen views thanks to our augmentations.
- ② Compared with recent state-of-the-art MipNeRF and other classical approaches, Aug-NeRF shows a clear

Table 2. Quantitative comparison of our Aug-NeRF against NeRF and other top-performer for novel view synthesis. Performance is reported on the test set of 360° NeRF-Sythetic dataset. \uparrow/\downarrow means that larger/smaller numbers denote better performance.

Scene “drums”	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Average \downarrow
SRN [52]	17.18	0.766	0.267	0.135
NV [21]	22.58	0.873	0.214	0.075
LLFF [52]	21.13	0.890	0.126	0.069
NeRF [29]	25.01	0.925	0.091	0.043
MipNeRF [3]	26.22	0.939	0.065	0.034
Aug-NeRF (Ours)	26.34	0.941	0.060	0.032
Scene “ship”	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Average \downarrow
SRN [52]	20.60	0.757	0.299	0.109
NV [21]	23.93	0.784	0.276	0.080
LLFF [52]	23.22	0.823	0.218	0.076
NeRF [29]	28.65	0.856	0.206	0.047
MipNeRF [3]	29.30	0.864	0.190	0.044
Aug-NeRF (Ours)	29.18	0.879	0.173	0.042
Scene “chair”	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Average \downarrow
SRN [52]	26.96	0.910	0.106	0.040
NV [21]	28.33	0.916	0.109	0.036
LLFF [52]	28.72	0.948	0.064	0.027
NeRF [29]	33.00	0.967	0.046	0.016
MipNeRF [3]	33.82	0.972	0.042	0.014
Aug-NeRF (Ours)	33.87	0.972	0.040	0.014

Table 3. Quantitative comparison of Aug-NeRF against NeRF and MipNeRF for novel view synthesis. All models are trained with noisy data. Performance is reported on the *noise-free* test set.

“fern” + Gaussian Noise	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Average \downarrow
NeRF [29]	16.95	0.451	0.535	0.200
Aug-NeRF (Ours)	17.12	0.535	0.495	0.187
“fern” + Shot Noise	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Average \downarrow
NeRF [29]	15.75	0.231	0.755	0.260
Aug-NeRF (Ours)	17.00	0.495	0.485	0.190

advantage, especially in terms of PSNR. In some cases, MipNeRF has a slightly higher SSIM; but Aug-NeRF is able to outperform it in most cases.

- ③ Aug-NeRF achieves superior performance in representing fine geometry, as shown in Fig. 3 such as Fern’s and Orchid’s leaves, the skeleton ribs, and railing in T-rex. Both NeRF and MipNeRF reconstruct the low-frequency geometry and color variation, but fail to generate high-quality fine details (see zoom-in).

Depth and Geometry visualization. The learned depth maps and fitted 3D geometries from NeRFs are provided in and 4 and Fig. 5, respectively. The 3D shapes (Fig. 5) are synthesized by MarchingCube algorithms [22]. We observe that vanilla NeRF suffers from a serrated surface (which overwhelms the fine details), while traditional TV and Laplacian regularizations tend to excessively smoothen the results. Aug-NeRF reduces noises and improves surface smoothness, in a detail- and geometry-preserving manner.

Table 4. Quantitative ablation study of our Aug-NeRF. Input, feature, and output augmentations. denote our proposed coordinate, feature, and pre-rendering output augmentations respectively.

Scene “fern”	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Average \downarrow
NeRF [29]	25.17	0.792	0.280	0.073
+ ℓ_1 Reg.	25.15	0.750	0.285	0.076
+ Lap. Reg.	24.89	0.670	0.305	0.083
+ TV Reg.	26.05	0.806	0.217	0.062
+ Random Aug.	25.28	0.796	0.224	0.067
+ Input Aug.	25.30	0.797	0.251	0.069
+ Feature Aug.	25.39	0.787	0.243	0.069
+ Feature & (Pre-) Output Aug.	26.32	0.810	0.199	0.059
+ Tri-level Random Noise	25.36	0.802	0.205	0.064
Aug-NeRF (Ours)	26.51	0.830	0.168	0.054

Superior synthesis when trained on noisy data. As an extra study, we examine Aug-NeRF under supervision images with additive noise corruptions. From Tab. 3 and Fig. 6, compared to the vanilla NeRF, Aug-NeRF shows consistent 6.5% \sim 26.9% average error reductions for both Gaussian and Shot noises, while it substantially improves the visual quality of constructed test views (e.g., much fewer noises in the “fern”). We regard it as an additional bonus from enforcing smooth geometry in NeRF training.

5.3. Ablation Study

Multi-level v.s. single-level augmentation. To compare the effects of robust augmentations at different levels, we conduct step-wise evaluation as: (i) NeRF, (ii) NeRF + Feature Aug., (iii) NeRF + Feature & Output Aug., (iv) NeRF + Feature & Output & Input coordinates Aug., which is our complete Aug-NeRF. Tab. 4 shows that applying robust augmentation to each level brings extra and complementary generalization gains, among which augmenting the pre-rendering output level makes the biggest difference.

Worst-case v.s. random perturbations. One straightforward baseline for Aug-NeRF is to just use random data augmentation. Particularly, we employ random Gaussian noises to both intermediate features and pre-rendering outputs of NeRF¹. As in Tab. 4, Random Aug. obtains moderate performance boosts for all metrics, but are clearly less obvious than our worst-case perturbations.

Effects of augmentation strength and location. The accuracy gains from Aug-NeRF are largely determined by the strength and location of crafted worst-case perturbations. A comprehensive investigation on three levels of augmentations, i.e., input coordinate, intermediates features, and pre-rendering output, are presented in Fig. 7. When studying one of the factors, we stick to the best configuration for the rest factors. Fig. 7 reveals that: First, NeRF gains the most from {coordinate, features, pre-rendering output} augmentations with {PGD-3, PGD-1, PGD-1} and step size $\{10^{-2}, 10^{-3}, 10^{-5}\}$; Second, applying generated perturbations to the middle layer of NeRF’s MLP contributes the most significantly;

¹The vanilla NeRF has already included random noise in coordinates.

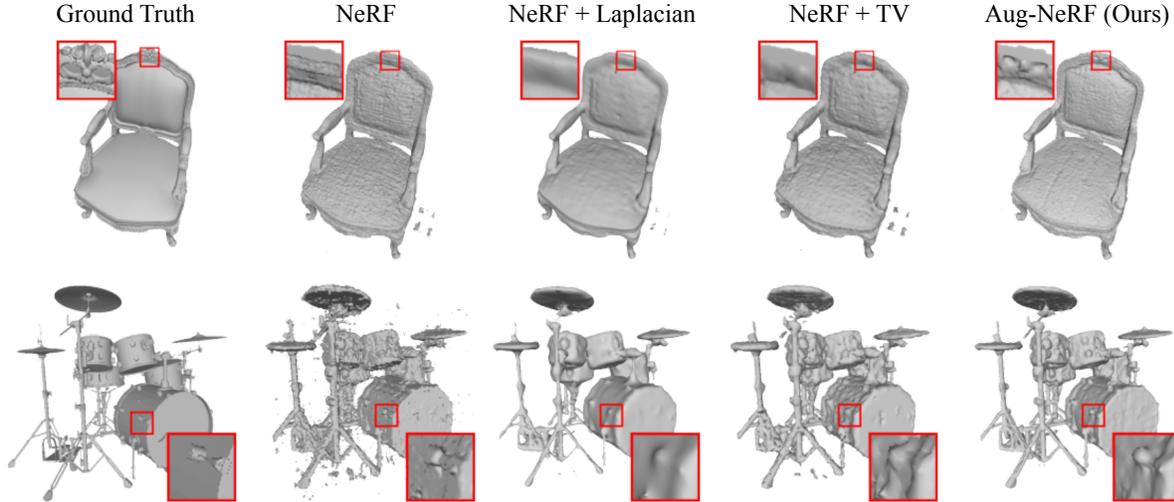


Figure 5. Comparisons of the fitted geometry by NeRF, NeRF with explicit Laplacian and TV regularizers, and our Aug-NeRF. The local zoom-in is placed in the red box.

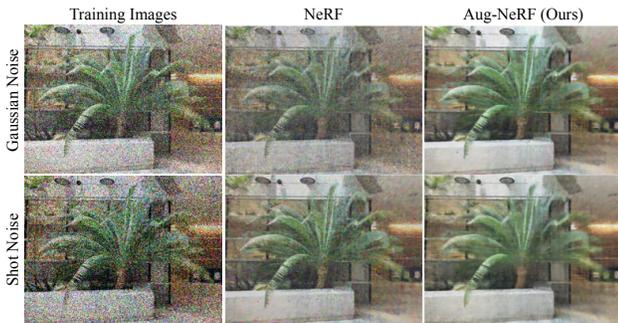


Figure 6. AugNeRF yields superior synthesis results when trained with noisy image supervision. Our investigated image corruptions include Gaussian and shot noises, following the standard in [13].

Third, too strong (e.g., PGD-10) worst-case perturbations may still deteriorate performance.

Comparison with explicit smooth regularizations. In contrast to our implicit smooth prior, there exists several explicit smooth regularizations which can be directly plugged into the NeRF pipeline, like ① ℓ_1 sparsity Reg. $R_{\ell_1}(\Theta) = \int |\sigma_{\Theta}(\mathbf{u})| d\mathbf{u}$; ② Laplacian Reg. $R_{\text{Lap}}(\Theta) = \int |\Delta \sigma_{\Theta}(\mathbf{u})| d\mathbf{u}$; ③ Total Variation (TV) Reg. $R_{\text{TV}}(\Theta) = \int \|\nabla \sigma_{\Theta}(\mathbf{u})\|_2 d\mathbf{u}$. As demonstrated in Tab. 4, although hyperparameters are carefully tuned by a grid search, both ℓ_1 and Laplacian regularizers degrade the performance, as such explicit constraints are often too aggressive and limit the representation flexibility of NeRF. The TV Reg. can lead to positive gains but still largely lags behind our proposals.

6. Conclusion and Broad Impact

In this paper, we have presented Aug-NeRF that addresses the inherent non-smooth geometries of NeRF. Specifically, based on solid physical grounds, Aug-NeRF seamlessly injects worst-case perturbations into three levels of the NeRF

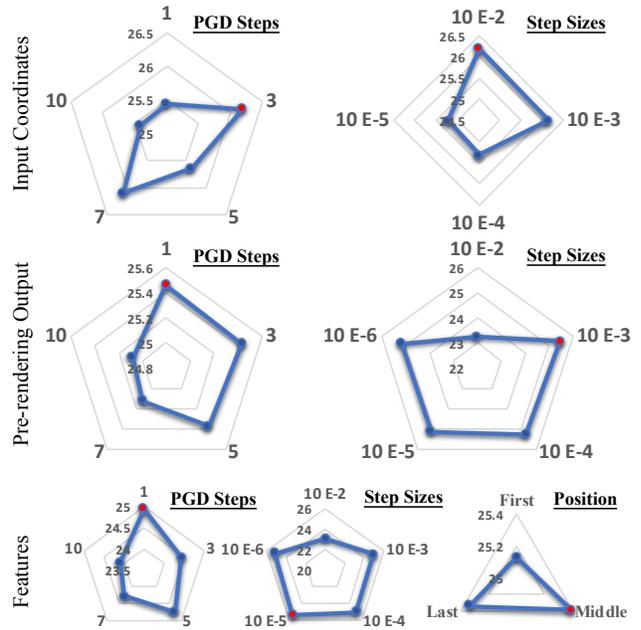


Figure 7. Ablations on the strength and location of all three augmentations. Results are on the test set of LLFF. The PGD step is the number of iterations for generating worst-case perturbations. Step size α controls the strength of crafted perturbations. Position means which layer the intermediate feature augmentation is injected into NeRF’s MLP. Red indicates the top performance.

pipeline, leading to substantially improved geometry continuity and generalization ability. Extensive quantitative and qualitative results across diverse scenes validate the effectiveness of our proposals. Moreover, the implicit smooth prior induced by triple-level augmentation enables NeRF to recover scenes from noisy supervision images. One limitation is that we only study additive noises (e.g., Gaussian) for corrupted images. We will extend the investigation to other complicated corruptions.

References

- [1] Kara-Ali Aliev, Artem Sevastopolsky, Maria Kolos, Dmitry Ulyanov, and Victor Lempitsky. Neural point-based graphics. In *European Conference on Computer Vision (ECCV)*, 2020. 2
- [2] Relja Arandjelović and Andrew Zisserman. Nerf in detail: Learning to sample for view synthesis. *arXiv preprint arXiv:2106.05264*, 2021. 4
- [3] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. 2, 6, 7
- [4] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. 2
- [5] Tianlong Chen, Yu Cheng, Zhe Gan, Jianfeng Wang, Lijuan Wang, Zhangyang Wang, and Jingjing Liu. Adversarial feature augmentation and normalization for visual recognition. *arXiv preprint arXiv:2103.12171*, 2021. 2, 4, 5
- [6] Julian Chibane, Thiemo Alldieck, and Gerard Pons-Moll. Implicit functions in feature space for 3d shape reconstruction and completion. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [7] Daniel Cremers, Mikael Rousson, and Rachid Deriche. A review of statistical approaches to level set segmentation: integrating color, texture, motion and shape. *International Journal of Computer Vision (IJCV)*, 2007. 4
- [8] Farzan Farnia, Jesse M Zhang, and David Tse. Generalizable adversarial training via spectral normalization. In *International Conference on Learning Representations (ICLR)*, 2019. 2
- [9] Zhe Gan, Yen-Chun Chen, Linjie Li, Chen Zhu, Yu Cheng, and Jingjing Liu. Large-scale adversarial training for vision-and-language representation learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2, 4, 5
- [10] Chen Gao, Yichang Shih, Wei-Sheng Lai, Chia-Kai Liang, and Jia-Bin Huang. Portrait neural radiance fields from a single image. *arXiv preprint arXiv:2012.05903*, 2020. 2
- [11] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations (ICLR)*, 2015. 2
- [12] Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuozhuo Dai, Feitong Tan, and Ping Tan. Cascade cost volume for high-resolution multi-view stereo and stereo matching. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [13] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. In *International Conference on Learning Representations (ICLR)*, 2019. 8
- [14] Dan Hendrycks, Norman Mu, Ekin Dogus Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. Augmix: A simple data processing method to improve robustness and uncertainty. In *International Conference on Learning Representations (ICLR)*, 2019. 2
- [15] Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. In *International Conference on Learning Representations (ICLR)*, 2020. 4
- [16] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial examples in the physical world. In *International Conference on Learning Representations Workshop (ICLRW)*, 2017. 2
- [17] Uday Kusupati, Shuo Cheng, Rui Chen, and Hao Su. Normal assisted stereo depth estimation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 1
- [18] Kiriakos N Kutulakos and Steven M Seitz. A theory of shape by space carving. *International Journal of Computer Vision (IJCV)*, 2000. 2
- [19] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss landscape of neural nets. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. 4
- [20] Shichen Liu, Shunsuke Saito, Weikai Chen, and Hao Li. Learning to infer implicit surfaces without 3d supervision. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019. 2
- [21] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. In *ACM Conference and Exhibition on Computer Graphics and Interactive Techniques (SIGGRAPH)*, 2019. 2, 6, 7
- [22] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. In *ACM Conference and Exhibition on Computer Graphics and Interactive Techniques (SIGGRAPH)*, 1987. 7, A14
- [23] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations (ICLR)*, 2018. 2
- [24] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [25] Nelson Max. Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics (TVCG)*, 1995. 3
- [26] Quan Meng, Anpei Chen, Haimin Luo, Minye Wu, Hao Su, Lan Xu, Xuming He, and Jingyi Yu. Gnerf: Gan-based neural radiance field without posed camera. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. 2
- [27] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [28] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and

- Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 2019. 2, 6
- [29] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European Conference on Computer Vision (ECCV)*, 2020. 1, 2, 4, 5, 6, 7, A13
- [30] Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, and Shin Ishii. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2018. 4
- [31] Preetum Nakkiran. Adversarial robustness may be at odds with simplicity. *arXiv preprint arXiv:1901.00532*, 2019. 2
- [32] Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nati Srebro. Exploring generalization in deep learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. 4
- [33] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [34] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5589–5599, 2021. 2
- [35] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [36] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. In *ACM Conference and Exhibition on Computer Graphics and Interactive Techniques in Asia (SIGGRAPH Asia)*, 2021. 2
- [37] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *Advances in Neural Information Processing Systems Workshop (NeurIPSW)*, 2017. A14
- [38] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *European Conference on Computer Vision (ECCV)*, 2020. 2
- [39] Eric Penner and Li Zhang. Soft 3d reconstruction for view synthesis. *ACM Transactions on Graphics (TOG)*, 2017. 2
- [40] Aditi Raghunathan, Sang Michael Xie, Fanny Yang, John C Duchi, and Percy Liang. Adversarial training can hurt generalization. In *International Conference on Machine Learning Workshop (ICMLW)*, 2019. 2
- [41] Amit Raj, Michael Zollhofer, Tomas Simon, Jason Saragih, Shunsuke Saito, James Hays, and Stephen Lombardi. Pva: Pixel-aligned volumetric avatars. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [42] Sylvestre-Alvise Rebuffi, Sven Gowal, Dan A Calian, Florian Stimberg, Olivia Wiles, and Timothy Mann. Fixing data augmentation to improve adversarial robustness. *arXiv preprint arXiv:2103.01946*, 2021. 2
- [43] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. 2
- [44] Konstantinos Rematas, Ricardo Martin-Brualla, and Vittorio Ferrari. Sharf: Shape-conditioned radiance fields from a single view. In *International Conference on Machine Learning (ICML)*, 2021. 2
- [45] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *European Conference on Computer Vision (ECCV)*, 2020. 2
- [46] Gernot Riegler and Vladlen Koltun. Stable view synthesis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [47] Andrea Romanoni and Matteo Matteucci. Tapa-mvs: Textureless-aware patchmatch multi-view stereo. In *IEEE International Conference on Computer Vision (ICCV)*, 2019. 1
- [48] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *IEEE International Conference on Computer Vision (ICCV)*, 2019. 2
- [49] Ludwig Schmidt, Shibani Santurkar, Dimitris Tsipras, Kunal Talwar, and Aleksander Madry. Adversarially robust generalization requires more data. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018. 2
- [50] Steven M Seitz and Charles R Dyer. Photorealistic scene reconstruction by voxel coloring. *International Journal of Computer Vision (IJCV)*, 1999. 2
- [51] Rakesh Shrestha, Zhiwen Fan, Qingkun Su, Zuozhuo Dai, Siyu Zhu, and Ping Tan. Meshmvs: Multi-view stereo guided mesh reconstruction. In *International Conference on 3D Vision (3DV)*, 2021. 2
- [52] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2, 6, 7
- [53] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhofer. Deepvoxels: Learning persistent 3d feature embeddings. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [54] Olga Sorkine, Daniel Cohen-Or, Yaron Lipman, Marc Alexa, Christian Rössl, and H-P Seidel. Laplacian surface editing. In *ACM SIGGRAPH Symposium on Geometry Processing*, 2004. 4
- [55] Pratul P Srinivasan, Ben Mildenhall, Matthew Tancik, Jonathan T Barron, Richard Tucker, and Noah Snavely. Lighthouse: Predicting lighting volumes for spatially-coherent illumination. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

- [56] Pratul P Srinivasan, Richard Tucker, Jonathan T Barron, Ravi Ramamoorthi, Ren Ng, and Noah Snavely. Pushing the boundaries of view extrapolation with multiplane images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [57] David Stutz, Matthias Hein, and Bernt Schiele. Disentangling adversarial robustness and generalization. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [58] David Stutz, Matthias Hein, and Bernt Schiele. Relating adversarially robust generalization to flat minima. *IEEE International Conference on Computer Vision (ICCV)*, 2021. 4
- [59] Ke Sun, Zhanxing Zhu, and Zhouchen Lin. Towards understanding adversarial examples systematically: Exploring data size, task and model factors. *arXiv preprint arXiv:1902.11019*, 2019. 2
- [60] Richard Szeliski and Polina Golland. Stereo matching with transparency and matting. In *International Conference on Computer Vision (ICCV)*, 1998. 2
- [61] Matthew Tancik, Ben Mildenhall, Terrance Wang, Divi Schmidt, Pratul P Srinivasan, Jonathan T Barron, and Ren Ng. Learned initializations for optimizing coordinate-based neural representations. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [62] Matthew Tancik, Pratul P Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2, 3, 4
- [63] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *ACM Transactions on Graphics (TOG)*, 2019. 2
- [64] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. In *International Conference on Learning Representations (ICLR)*, 2019. 2
- [65] Dilin Wang, Chengyue Gong, and Qiang Liu. Improving neural language modeling via adversarial training. In *International Conference on Machine Learning (ICML)*, 2019. 2, 4, 5
- [66] Haotao Wang, Chaowei Xiao, Jean Kossaifi, Zhiding Yu, Anima Anandkumar, and Zhangyang Wang. Augmax: Adversarial composition of random augmentations for robust training. *Advances in Neural Information Processing Systems (NeurIPS)*, 34, 2021. 2
- [67] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689*, 2021. 2
- [68] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [69] Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing (TIP)*, 2004. 6
- [70] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. Nerf-: Neural radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064*, 2021. 2
- [71] Colin Wei and Tengyu Ma. Improved sample complexities for deep networks and robust classification via an all-layer margin. In *International Conference on Learning Representations (ICLR)*, 2020. 2, 4, 5
- [72] Yi Wei, Shaohui Liu, Yongming Rao, Wang Zhao, Jiwen Lu, and Jie Zhou. Nerfingmvs: Guided optimization of neural radiance fields for indoor multi-view stereo. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. 2
- [73] Dongxian Wu, Shu-Tao Xia, and Yisen Wang. Adversarial weight perturbation helps robust generalization. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2
- [74] Minye Wu, Yuehao Wang, Qiang Hu, and Jingyi Yu. Multi-view neural human rendering. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [75] Cihang Xie, Mingxing Tan, Boqing Gong, Jiang Wang, Alan L Yuille, and Quoc V Le. Adversarial examples improve image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 4, 5
- [76] Guandao Yang, Serge Belongie, Bharath Hariharan, and Vladlen Koltun. Geometry processing with neural fields. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 4
- [77] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2018. 2
- [78] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems*, 34:4805–4815, 2021. 2
- [79] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Ronen Basri, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2
- [80] Lin Yen-Chen, Pete Florence, Jonathan T Barron, Alberto Rodriguez, Phillip Isola, and Tsung-Yi Lin. Inerf: Inverting neural radiance fields for pose estimation. In *IEEE International Conference on Intelligent Robots and Systems (IROS)*, 2020. 2
- [81] Wang Yifan, Shihao Wu, Cengiz Oztireli, and Olga Sorkine-Hornung. Iso-points: Optimizing neural implicit surfaces with hybrid representations. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 4
- [82] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. 2

- [83] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. [2](#)
- [84] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P Xing, Laurent El Ghaoui, and Michael I Jordan. Theoretically principled trade-off between robustness and accuracy. In *International Conference on Machine Learning (ICML)*, 2019. [2](#)
- [85] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. [2](#)
- [86] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. [6](#)
- [87] Xiaofeng Zhang, Zhangyang Wang, Dong Liu, and Qing Ling. Dada: Deep adversarial data augmentation for extremely low data regime classification. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019. [2](#)
- [88] Chao Zhou, Hong Zhang, Xiaoyong Shen, and Jiaya Jia. Unsupervised learning of stereo matching. In *IEEE International Conference on Computer Vision (ICCV)*, 2017. [1](#), [4](#)
- [89] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. In *ACM Conference and Exhibition on Computer Graphics and Interactive Techniques (SIGGRAPH)*, 2018. [2](#)
- [90] Chen Zhu, Yu Cheng, Zhe Gan, Siqi Sun, Tom Goldstein, and Jingjing Liu. Freelj: Enhanced adversarial training for natural language understanding. In *International Conference on Learning Representations (ICLR)*, 2020. [2](#), [4](#), [5](#)

A1. More Technical Details

We summarize the detailed procedures of Aug-NeRF in the Algorithm 1.

Algorithm 1 The training pipeline of Aug-NeRF. For simplicity, we assume batch size is 1.

Initialize: Training view images $\mathcal{I} = \{I_i \in \mathbb{R}^M\}_{i=1}^N$ and their associated camera poses $\mathcal{P} = \{\phi_i \in \mathbb{R}^{3 \times 4}\}_{i=1}^N$. Define neural radiance field $F_{\Theta}(\mathbf{p}, \theta) = (g \circ f(\mathbf{p}, \theta), h \circ f(\mathbf{p})) : (\mathbf{p}, \theta) \mapsto (c, \sigma)$ as in Sec. 4.2.

- 1: Cast rays for each pixel in each I_i via inverse projection with respect to ϕ_i , and obtain a set of rays $\mathcal{R} = \{(\mathbf{o}_i, \mathbf{d}_i, \theta_i, \widehat{C}_i)\}_{i=1}^{NM}$.
- 2: **while** until convergence **do**
- 3: Randomly pick a ray $(\mathbf{o}_i, \mathbf{d}_i, \theta_i, \widehat{C}_i) \in \mathcal{R}$
- 4: Generate adversarial perturbations $\delta_t, \delta_{xyz}, \delta_\theta, \delta_f, \delta_c, \delta_\sigma$ by solving Eqn. 4 using PGD (Eqn. 5) within the corresponding search space in Sec. 4.2.
- 5: # Sample points along rays.
- 6: (Coarse) Sample $K/2$ depth intervals t_k along the rays uniformly
- 7: (Fine) Sample $K/2$ depth intervals t_k via proportional to coarse sampled densities [29].
- 8: **for** $k \in \{1, \dots, K\}$ **do**
- 9: $t_k^\dagger = t_k + \delta_{t,k}, \theta_k^\dagger = \theta_i + \delta_\theta.$
- 10: $\mathbf{p}_k = \mathbf{o}_i + t_k \mathbf{d}_i, \mathbf{p}_k^\dagger = \mathbf{o}_i + t_k^\dagger \mathbf{d}_i + \delta_{xyz}$
- 11: $(\mathbf{c}_k, \sigma_k) = g \circ f(\mathbf{p}_k, \theta_i), h \circ f(\mathbf{p}_k)$
- 12: $\mathbf{c}_k^\dagger = g(f(\mathbf{p}_k^\dagger) + \delta_f, \theta_k^\dagger) + \delta_c$
- 13: $\sigma_k^\dagger = h(f(\mathbf{p}_k^\dagger + \delta_f)) + \delta_\sigma$
- 14: **end for**
- 15: # Volumetric rendering.
- 16: $(\Delta t_k, \Delta t_k^\dagger) = t_k - t_{k-1}, t_k^\dagger - t_{k-1}^\dagger$
- 17: $T(k) = \exp\left(-\sum_{l=1}^{k-1} \sigma_l \Delta t_l\right)$
- 18: $T^\dagger(k) = \exp\left(-\sum_{l=1}^{k-1} \sigma_l^\dagger \Delta t_l^\dagger\right)$
- 19: $\mathbf{C}_i = \sum_{k=1}^K T(k)(1 - \exp(-\sigma_k \Delta t_k)) \mathbf{c}_k$
- 20: $\mathbf{C}_i^\dagger = \sum_{k=1}^K T^\dagger(k)(1 - \exp(-\sigma_k^\dagger \Delta t_k^\dagger)) \mathbf{c}_k^\dagger$
- 21: # Train network.
- 22: $\mathcal{L} = \|\mathbf{C}_i - \widehat{C}_i\|_2^2 + \lambda \|\mathbf{C}_i^\dagger - \widehat{C}_i\|_2^2$
- 23: Update network parameter Θ via $\nabla_{\Theta} \mathcal{L}$.
- 24: **end while**

A2. More Experiment Results

Qualitative results on NeRF-Synthetic 360° dataset. We present the constructed test views in Fig. A8 and the learned depth maps in Fig. A9. As shown in Fig. A8, we find that the vanilla NeRF fails to capture the fine-grained details of objects, such as the “ship net”, while Aug-NeRF demonstrates substantially improved visual qualities.

In the meantime, from the depth maps in Fig. A9, NeRF

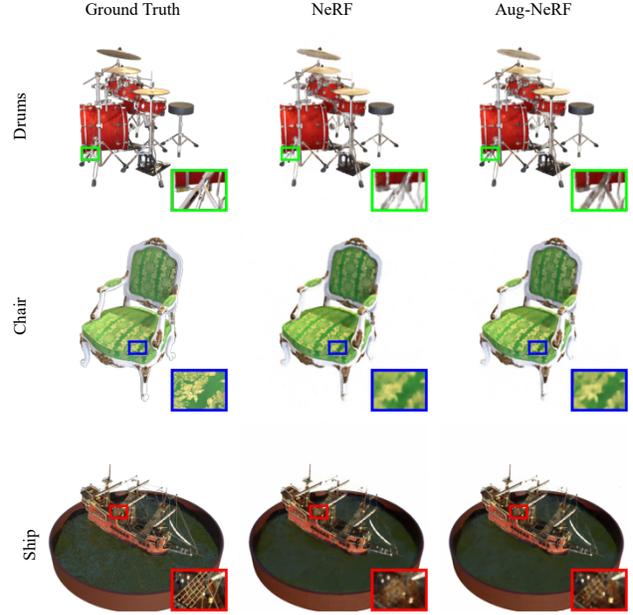


Figure A8. Comparisons on test-set views for scenes from the NeRF-Synthetic 360° dataset generated with a physics-based renderer [29].

baseline suffers from severe noises. On the contrary, Aug-NeRF enjoys much more smooth depth maps, which suggests that our proposed triple-level robust augmentations indeed enhance the NeRF’s continuity and generate smooth geometry representations.

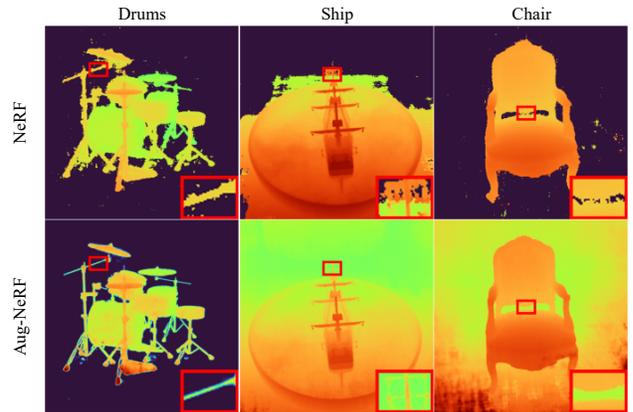


Figure A9. Comparisons of learned depth maps from NeRF and our Aug-NeRF on scenes from NeRF-Synthetic 360° dataset.

Benefits in overfitting vs underfitting cases. We take the scene “chair” as an example, and investigate three combinations of different model sizes and data scales as below Tab. A5: (i) big NeRF (512) on small images ($\frac{1}{2}$ Res.); (ii) normal NeRF (256) on small images ($\frac{1}{2}$ Res.); (iii) small NeRF (128) on large images (Full Res.). We show that in all settings of overfitting / normal case / underfitting, our

proposed augmentations are consistently beneficial.

Table A5. Performance of Aug-NeRF on different backbone and data size combinations.

Setting	Model	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Average \downarrow
Big NeRF (512) + $\frac{1}{2}$ Res.	NeRF	33.56	0.968	0.043	0.015
	Aug-NeRF (Ours)	34.26	0.973	0.038	0.013
Big NeRF (256) + $\frac{1}{2}$ Res.	NeRF	33.27	0.968	0.045	0.016
	Aug-NeRF (Ours)	33.92	0.971	0.038	0.013
Small NeRF (128) + Full Res.	NeRF	33.00	0.967	0.046	0.016
	Aug-NeRF (Ours)	33.86	0.970	0.041	0.014

Geometry extraction. To obtain geometric visualization in Fig. 5, we first query the network F_Θ with a regular lattice defined over $[-1, 1]^3$, and export a discretized density field volume. The absolute voxel size is $2/512$. Then we employ marching cube algorithm [22] provided in UCSF Chimera² to extract the surface. We set the threshold to 25 and 1 for chairs and drums, respectively. The step size is chosen as 1. In order to numerically assess the quality of reconstructed geometries, we introduce Chamfer Distance (CD) to measure the difference between reconstructed geometries and ground-truth models:

$$d_{CD} = \frac{1}{|\mathcal{S}_1|} \sum_{\mathbf{x} \in \mathcal{S}_1} \min_{\mathbf{y} \in \mathcal{S}_2} \|\mathbf{x} - \mathbf{y}\|_2 + \frac{1}{|\mathcal{S}_2|} \sum_{\mathbf{x} \in \mathcal{S}_2} \min_{\mathbf{y} \in \mathcal{S}_1} \|\mathbf{x} - \mathbf{y}\|_2,$$

where \mathcal{S}_1 and \mathcal{S}_2 are point sets sampled from the extracted surfaces and ground-truth models, respectively. On scene chair, our AugNeRF achieves 1.04×10^{-2} CD which is 29.25% lower than vanilla NeRF (1.47×10^{-2}).

Different types of noise and inaccurate camera poses.

As shown in Tab. 3 and Fig. 6, we experiment on two kinds of corruptions, i.e., Gaussian and Shot noises. In this paragraph, we add extra results of training with Pepper noise and inaccurate camera poses are collected in below Tab. A6. The results consistently demonstrate the superiority of Aug-NeRF. We note that our main goal is to endow NeRF with smoothness-aware geometry reconstruction, enhanced generalization to synthesizing unseen views, while the improved tolerance of noisy supervisions is a by-product bonus.

Table A6. Additional results of Aug-NeRF trained on images corrupted by pepper noise and inaccurate camera poses.

Noise Type	“ferm”	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	Average \downarrow
Pepper Noise	NeRF	19.01	0.401	0.546	0.174
	Aug-NeRF (Ours)	19.96	0.568	0.403	0.138
Inaccurate Pose	NeRF	12.31	0.253	0.725	0.333
	Aug-NeRF (Ours)	13.54	0.365	0.811	0.306

Implementation of explicit regularization. We investigate three types of explicit regularizations: ℓ_1 sparsity, total

variation (TV), and Laplacian. The TV regularization is defined as:

$$R_{TV}(\Theta) = \int_{[-1,1]^3} |\nabla_{\mathbf{x}} \sigma_\Theta(\mathbf{x})| d\mathbf{x},$$

where σ_Θ denotes the density branch of the function F_Θ . However, evaluating this integral is implausible. Instead, we discretize the integral interval into regular grids and conducting quadrature rule for estimating TV regularization:

$$R_{TV}(\Theta) = \sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^N |\nabla_{\mathbf{x}} \sigma_\Theta(\delta i, \delta j, \delta k)| \delta^3,$$

where we utilize auto-differentiation provided in PyTorch Library [37] to calculate $\nabla_{\mathbf{x}} \sigma_\Theta$. Similarly, we can approximate ℓ_1 sparsity and Laplacian regularization by:

$$\begin{aligned} R_{\ell_1}(\Theta) &= \int_{[-1,1]^3} |\sigma_\Theta(\mathbf{x})| d\mathbf{x} \\ &\approx \sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^N |\sigma_\Theta(\delta i, \delta j, \delta k)| \delta^3 \\ R_{Lap}(\Theta) &= \int_{[-1,1]^3} |\Delta_{\mathbf{x}} \sigma_\Theta(\mathbf{x})| d\mathbf{x} \\ &\approx \sum_{i=1}^N \sum_{j=1}^N \sum_{k=1}^N |\Delta_{\mathbf{x}} \sigma_\Theta(\delta i, \delta j, \delta k)| \delta^3 \end{aligned}$$

where $\Delta = \text{div} \cdot \nabla$ denotes the Laplacian operator.

²<https://www.cgl.ucsf.edu/chimera/>