

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: $t_k^{\dagger} = t_k + \delta_{t,k}, \theta_i^{\dagger} = \theta_i + \delta_{\theta}$.
- 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: $(c_k, \sigma_k) = g \circ f(\mathbf{p}_k, \theta_i), h \circ f(\mathbf{p}_k)$
- 12: $c_k^{\dagger} = g(f(\mathbf{p}_k^{\dagger}) + \delta_f, \theta_i^{\dagger}) + \delta_c$
- 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: $C_i = \sum_{k=1}^K T(k)(1 - \exp(-\sigma_k \Delta t_k))c_k$
- 20: $C_i^{\dagger} = \sum_{k=1}^K T^{\dagger}(k)(1 - \exp(-\sigma_k^{\dagger} \Delta t_k^{\dagger}))c_k^{\dagger}$
- 21: # Train network.
- 22: $\mathcal{L} = \|C_i - \widehat{C}_i\|_2^2 + \lambda \|C_i^{\dagger} - \widehat{C}_i\|_2^2$
- 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

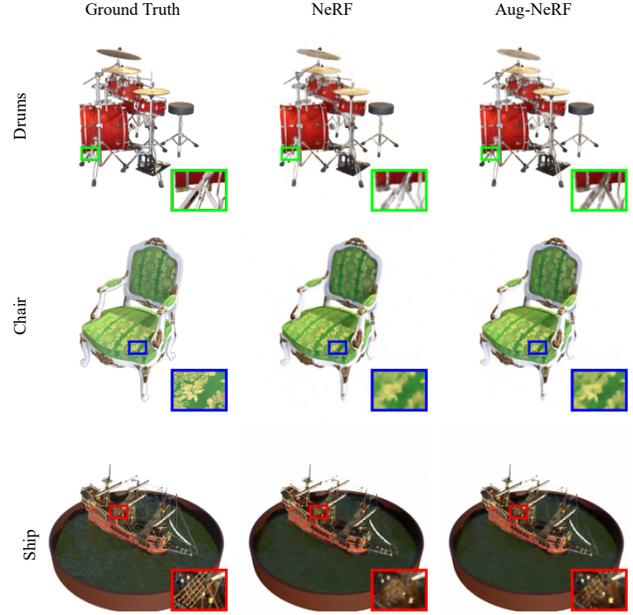


Figure A8. Comparisons on test-set views for scenes from the NeRF-Synthetic 360° dataset generated with a physics-based renderer [29].

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.

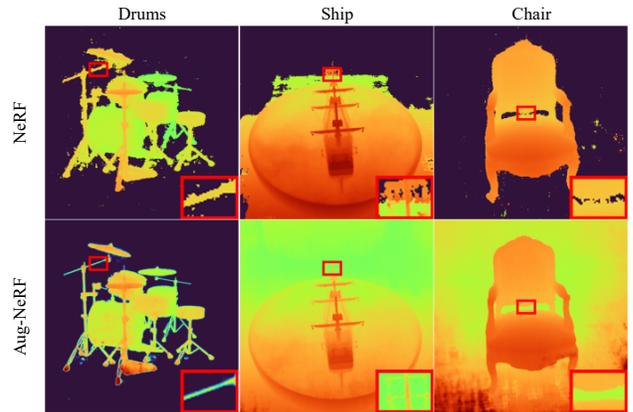


Figure A9. Comparisons of learned depth maps from NeRF and our Aug-NeRF on scenes from NeRF-Synthetic 360° dataset.

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 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/>