Advancing Super-Resolution in Neural Radiance Fields via Variational Diffusion Strategies

Shrey Vishen Monta Vista High School Cupertino, California

shrey.vishen@gmail.com

Chinmay Bharathulwar John P Stevens High School Edison, New Jersey

cbharathulwar3@gmail.com

Jatin Sarabu Bellarmine College Preparatory Santa Clara, California

jatinsarabu13@gmail.com

Rithwick Lakshmanan Pleasant Valley High School Bettendorf, Iowa

rithwick.laks@gmail.com

Saurav Kumar UI Urbana-Champaign Urbana, Illinois

sauravk4@illinois.edu

Vishnu Srinivas Foothill High School Pleasanton, California

vishnuvishaks@gmail.com

Abstract

We present a novel method for diffusion-guided frameworks for view-consistent super-resolution (SR) in neural rendering. Our approach leverages existing 2D SR models in conjunction with advanced techniques such as Variational Score Distilling (VSD) and a LoRA finetuning helper, with spatial training to significantly boost the quality and consistency of upscaled 2D images compared to the previous methods in the literature, such as Renoised Score Distillation (RSD) proposed in DiSR-NeRF [7], or SDS proposed in DreamFusion. The VSD score facilitates precise fine-tuning of SR models, resulting in high-quality, view-consistent images. To address the common challenge of inconsistencies among independent SR 2D images, we integrate Iterative 3D Synchronization (I3DS) from the DiSR-NeRF framework [5]. Quantitative benchmarks and qualitative results on the LLFF dataset demonstrate the superior performance of our system compared to existing methods such as DiSR-NeRF. All our code is available at https://github.com/ shreyvish5678/SR-NeRF-with-Variational-Diffusion-Strategies

1. Introduction

Neural Radiance Fields (NeRFs) have revolutionized 3D scene rendering from 2D images, significantly impacting applications such as 3D reconstruction and virtual reality [10]. NeRFs use continuous volumetric functions optimized by neural networks to synthesize high-fidelity views [19] [30] [21]. However, scaling NeRF for superresolution, maintaining view consistency, and managing

high-dimensional data remains challenging [31] [32] [20] [27]. Somewhat recent advancements such as Mip-NeRF 360 [2] and TensoRF have addressed some of these challenges with NeRF, especially with unbounded scene rendering and tensor decomposition. Our work focuses specifically on enhancing super-resolution capabilities within the NeRF framework.

Previous enhancements to NeRF, like Score Distillation Sampling (SDS), have struggled with issues such as oversmoothing and computational inefficiency, limiting their ability to capture fine details. To address these problems, we introduce Variational Score Distillation (VSD), which models 3D scene parameters as probabilistic distributions rather than fixed values. This approach improves scene representation by leveraging diffusion models [4] [13] and incorporates low-rank adaptation (LoRA) for efficient fine-tuning of pre-trained models.

Our extensive experiments show that VSD significantly outperforms SDS and RSD in generating detailed and photorealistic NeRFs, enhancing visual quality and computational efficiency. Results from datasets like LLFF confirm these improvements. Additionally, we provide a detailed ablation study on our system's components, including LoRA-based fine-tuning [6] and hierarchical sampling strategies [22].

Our approach advances NeRF-based rendering, offering a new standard for high-resolution 3D scene generation with broad applications in entertainment, gaming, scientific visualization, and architectural design [1] [9] [15].

To summarize, the contributions of this paper are as follows.

 Using Pre-trained Stable Diffusion weights for 3D NeRF Render Upscaling

- Utilizing Low-rank adaptation with Mixed Precision training for steering outputs
- Using a version of VSD loss and I3DS for training and upscaling the NeRF outputs

2. Methodolgy

2.1. Pre-Requisites

2.1.1 Latent Encoding and Residual Learning

Latent encoding starts by extracting 2D projections from a lower-resolution Neural Radiance Field (NeRF). These 2D views are then passed through an encoder, transforming them into latent space representations that capture key features in a compressed form [17].

To improve image quality, we introduce learnable residual latents—additional vectors added to the original encoding. During training, these residuals are adjusted to correct errors or enhance specific features, refining the latent vectors over time. This iterative process, where residuals update with the model, boosts both image quality and consistency. The combined latents look like this:

$$x_0' = x_0 + h_\theta \tag{1}$$

2.1.2 Forward Diffusion Process

Once the latent vectors, including the residual latents, are generated, they undergo a forward diffusion process [18]. Noise is gradually added to the latent vectors over timesteps, governed by the equation:

$$x_t = \sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$
 (2)

This transforms the structured latent representation into a noisy state by the final timestep. The diffusion process simulates the challenge of reversing noise to recover the original image details, essential for the model to learn the data distribution.

In the subsequent denoising step, the noisy latents are refined to produce high-quality, view-consistent images [4] [13]. The overall goal of this process is to prepare the latent vectors for reconstruction, reversing the noise added during forward diffusion.

2.1.3 UNet-Based Prediction

In the model, the prediction of the final high-quality image relies heavily on the use of UNet architectures. The UNet models are employed to process the noisy latent vectors generated during the forward diffusion process and to reconstruct these into cleaner, more accurate representations of the original scene.

The first stage of latent prediction uses a frozen, pretrained UNet model trained on a large image denoising dataset. This makes it well-suited for processing the noisy latent vectors generated by the forward diffusion process. The pre-trained UNet takes as input the noisy latent vector \mathbf{x}_t , time embeddings \mathbf{t} , text embeddings \mathbf{y} with applied class labels \mathbf{c} , and the low-resolution image \mathbf{I}_{LR} , aiming to predict the noise added to the latent vector \mathbf{x}_0 also called ϵ .

Mathematically, this is expressed as:

$$\epsilon_{\phi} = f_{\phi}(\mathbf{x}_t, \mathbf{t}, \mathbf{y}^{\mathbf{c}}, \mathbf{I}_{LR}) \tag{3}$$

where f_{ϕ} is the function representing the pre-trained model, and ϕ represents its parameters.

This output serves as a reference for evaluating the finetuned model's performance.

The noisy latent vector \mathbf{x}_t is also processed by a Fine-Tuned UNet model, which includes learnable Low-rank Adaptation (LoRA) parameters [6]. These allow the model to efficiently adapt for high-quality image reconstruction. Unlike the frozen pre-trained model, the Fine-Tuned UNet is trained during this stage, with LoRA parameters enabling targeted adjustments to the network layers.

The Fine-Tuned UNet takes the same inputs: \mathbf{x}_t , time embeddings \mathbf{t} , text embeddings \mathbf{y} , and low-resolution image \mathbf{I}_{LR} , plus class labels \mathbf{c} to focus on task-specific features [8]. Its prediction is:

$$\epsilon_{\varphi} = f_{\varphi}(\mathbf{x}_t, \mathbf{t}, \mathbf{y}, \mathbf{I}_{LR}, \mathbf{c})$$
 (4)

where f_{φ} represents the function learned by the Fine-Tuned UNet with LoRA and φ represents its parameters.

To quantify the improvement in image quality provided by the Fine-Tuned UNet, we compute the Variational Score Distillation (VSD) loss, which measures the difference between the predictions of the pre-trained and fine-tuned models. We used an L1 variation of the VSD loss proposed in Prolific Dreamer [26], defined as:

$$\mathcal{L}_{\text{VSD}}(\theta) = \mathbb{E}_{t,\epsilon,\mathbf{c}} \left[\omega(t) \cdot \| \epsilon_{\phi} - \epsilon_{\varphi} \| \right] \tag{5}$$

where $\omega(t)$ is a weighting function that adjusts the importance of the loss based on the timestep t.

This loss is then backpropagated through the network, specifically targeting the LoRA parameters in the Fine-Tuned UNet. By minimizing this loss, the model learns to generate high-quality latent representations that closely match the ideal outputs while incorporating task-specific adjustments through the class labels. Then we obtrain the new residual latents:

$$\theta_{\text{new}} = \theta_{\text{old}} - \eta \nabla_{\theta} \mathcal{L}_{\text{VSD}}(\theta) \tag{6}$$

To add, every few steps, or every step, the LoRA parameters are fine-tuned as well, with the following equation:

$$\mathcal{L}_{\text{Diff}}(\theta) = \mathbb{E}_{t,\epsilon,\mathbf{c}} \left[\left(f_{\varphi}(\mathbf{x}_{t}, t, \mathbf{c}, \mathbf{I}_{LR}, \mathbf{y}) - \epsilon \right)^{2} \right]$$
 (7)

Where, \mathbf{x}_t was obtained using the forward diffusion process with \mathbf{x}_0 and ϵ . Then we can backpropagate this loss into the LoRA parameters, likewise:

$$\varphi_{\text{new}} = \varphi_{\text{old}} - \eta \nabla_{\varphi} \mathcal{L}_{\text{Diff}}(\theta)$$
 (8)

Now the refined latent vector $\hat{\mathbf{x}}_0$ from the Fine-Tuned UNet and the Pre-trained one is iteratively improved through the training process, leading to progressively higher-quality image outputs. The final prediction is not only a denoised version of the latent vector but also one that has been optimized for the specific rendering task at hand, thanks to the targeted adjustments made possible by the LoRA parameters.

This dual UNet approach, leveraging both pre-trained and fine-tuned models, ensures that the latent prediction process is both accurate and adaptable, ultimately contributing to the generation of high-resolution, view-consistent images that surpass the quality of those produced by existing methods [8] [25].

2.2. Low-rank Adaptation (LoRA) Fine-Tuning

Low-rank Adaptation (LoRA) is a powerful technique designed to fine-tune pre-trained neural networks efficiently, particularly in scenarios where extensive retraining or structural modifications are impractical [6]. In the context of Neural Radiance Fields (NeRFs), LoRA introduces trainable low-rank matrices into selected layers of the network, allowing for effective adaptation to new datasets or specific rendering tasks with minimal computational overhead.

The core idea behind LoRA is to augment the existing layers of a pre-trained NeRF model with additional trainable parameters that capture essential modifications without altering the original weights. Specifically, LoRA integrates two low-rank matrices, denoted as $A \in \mathbb{R}^{m \times r}$ and $B \in \mathbb{R}^{r \times n}$, where r is a rank that is significantly smaller than the dimensions of the original weight matrix $W \in \mathbb{R}^{m \times n}$. The low-rank matrices A and B are introduced in such a way that the original weight matrix W is modified as follows:

$$W' = W + AB \tag{9}$$

Here, W^\prime represents the new effective weight matrix after the application of LoRA. This adjustment allows the model to learn additional features or adapt to new data without having to retrain the entire network from scratch.

The low-rank matrices A and B are fine-tuned during the training process, while the original weights W remain fixed. This approach ensures that the adaptation process is both efficient and effective, targeting only the parameters necessary for the specific task. The gradients for the matrices A and B are computed using standard backpropagation

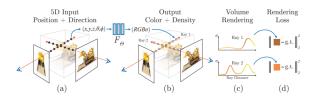


Figure 1. NeRF diagram

techniques, with the loss function \mathcal{L} defined for the specific rendering task, such as image enhancement, 3D reconstruction, or scene understanding. The gradients are given by:

$$\frac{\partial \mathcal{L}}{\partial A} = \frac{\partial \mathcal{L}}{\partial W'} B^T, \quad \frac{\partial \mathcal{L}}{\partial B} = A^T \frac{\partial \mathcal{L}}{\partial W'}$$
 (10)

These gradients guide the updates to A and B during the training process, allowing the NeRF model to fine-tune its performance on the new task without modifying the extensive pre-trained weights W.

2.3. Lower Resolution NeRF Generation

The first step in our pipeline involves constructing a low-resolution (LR) Neural Radiance Field (NeRF) from an LR image dataset. Using Nvidia's InstantNGP framework [12], we efficiently build this NeRF model. Afterward, the LR images are upsampled to match the resolution of the target super-resolution (SR) images, establishing a baseline for the iterative training process [14] [16] [3].

2.4. SR Training Process

The training process builds upon recent NeRF advancements, utilizing insights from a variety of volumetric rendering techniques [23] [29] and neural field applications [28]. The training loop operates iteratively until convergence, following a multi-step approach:

- 1. **Random Render Sampling**: A random LR image is rendered from the NeRF and used as input for encoding. This image is then taken and interpolated to 4x it's size, and then converted into a latent. Both the image and latent are passed. A text prompt is passed as well which describes the NeRF content along with from what orientation the image should look like, to condition the model for the desired output.
- 2. **Learnable Residual Latents**: Trainable residual latents are added to the latent encoding, for helping to refine image quality over time.
- 3. **Forward Diffusion**: The combined latents are passed through the forward diffusion process, to get the noisy latents, along with a given timestep embedding.
- 4. **Pre-trained UNet Prediction**: The noisy latent, along with text and time embeddings, is passed through a frozen, pre-trained UNet, which predicts a latent output used as a reference. The noisy latent is also processed by a

fine-tuned UNet model with Low-Rank Adaptation (LoRA) parameters to adapt the network efficiently. Class labels are added to guide the learning process.

- 5. **VSD Loss Calculation**: The Variational Score Distillation (VSD) loss, based on the L1 loss between the pretrained and fine-tuned UNet outputs, drives the optimization of the residual latents, by backpropagating this loss to the residual latents.
- 6. **Noise Differentiation Loss to LoRA**: Every few iterations, an auxiliary noise differentiation loss further refines the LoRA parameters. LoRA parameters are updated through backpropagation with this loss to improve task-specific fine-tuning.

This process repeats until the model converges, producing high-resolution, view-consistent NeRF outputs that outperform existing methods [11] [24], as we can see in Algorithm 1

Algorithm 1 VSD Super Resolution

```
1: Inputs: Latent x_0, text prompt embeddings y, timestep t, LR Image I_{lr}, class labels c, max timesteps M
```

```
2: Outputs: Latent residuals h_{\theta}
 3: Initialize h_{\theta}
 4: for S = [0, M] do
                 \epsilon \sim \mathcal{N}(0,1)
 5:
                x_0' = x_0 + h_\theta
 6:
                 x_t = \sqrt{a_t}x_0' + \sqrt{1 - a_t}\epsilon
 7:
                 \epsilon_{\phi} = f_{\phi}(x_t, t, y^c, I_{lr})
 8:
                 \epsilon_{\varphi} = f_{\varphi}(x_t, t, y, I_{lr}, c)
 9:
                \mathcal{L}_{VSD} = \mathbb{E}_{t,\epsilon,c} [w(t) \cdot || \epsilon_{\phi} - \epsilon_{\varphi} ||]
\theta \leftarrow \theta - \eta_1 \nabla_{\theta} \mathcal{L}_{VSD}
10:
11:
                \mathcal{L}_{Diff} = \mathbb{E}_{t,\epsilon,c} \left[ \left( f_{\varphi}(x_t, t, y, I_{lr}, c) - \epsilon \right)^2 \right]
12:
```

 $\varphi \leftarrow \varphi - \eta_2 \nabla_{\varphi} \mathcal{L}_{Diff}$

13:

14: **end for**

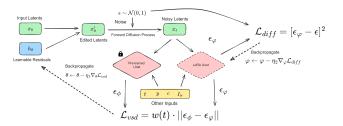


Figure 2. Refer to Algorithm 1

2.5. Iterative 3D Sync Process

The Iterative 3D Synchronization (I3DS) approach addresses limitations observed in initial experiments, where applying SDS directly on NeRF renders resulted in blurred details and convergence issues. We add I3DS to the pipeline

to improve this by decoupling the upscaling and NeRF synchronization processes into two alternating stages.

- **Upscaling Stage:** Starting with a low-resolution NeRF, images are rendered at 4x resolution. These images are then independently upscaled using RSD to add high-resolution details. However, initial upscaling may produce inconsistent details across views.
- Synchronization Stage: The upscaled images are used as inputs to update the NeRF model. During this stage, NeRF is trained using standard procedures to synchronize view-consistent details, correcting inconsistencies introduced during upscaling.

The synergy between upscaling and synchronization stages allows for increasingly detailed and view-consistent outputs over iterations. I3DS efficiently balances these stages to optimize for high-quality, consistent NeRF outputs while minimizing memory requirements and improving convergence times. This method demonstrates a reduction in optimization duration by 4x compared to previous approaches. We can see further about I3DS in Algorithm 2.

1: **Input:** LR NeRF ω_l , LR images I_{lr} , training poses P_{tr}

Algorithm 2 Iterative 3D Synchronization (I3DS)

```
2: Output: SR NeRF \omega_{sr}
 3: \omega \leftarrow \omega_l
 4: for S = [0, M] do
           Upscaling Stage
 5:
           x_0 \leftarrow \text{RenderImage}(\omega, P_{tr})
 7:
           x_0 \leftarrow \text{InterpolateX4}(x_0)
           z_0 \leftarrow \text{VaeEncode}(x_0)
 8:
           z_0' \leftarrow \text{VSD}(z_0, I_{lr})
 9:
           x_0' \leftarrow \text{VaeDecode}(z_0')
10:
           I_{tr} \leftarrow x_0'
11:
           Synchronization Stage
12:
           for sync_iter = [0, max\_sync\_iter] do
13:
                 (r_o, r_d, c_{tr}) \leftarrow \text{SampleRays}(I_{tr}, P_{tr})
14:
                 c' \leftarrow \text{RenderRays}(r_o, r_d)
15:
                 Perform gradient descent on \nabla_{\omega} \|c' - c_{tr}\|
16:
17:
                 \omega_{old} \leftarrow \omega
           end for
18:
19: end for
20: Return \omega_{sr} \leftarrow \omega
```

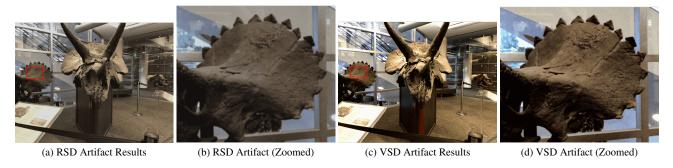


Figure 4. Comparison of RSD and VSD on Museum Artifact

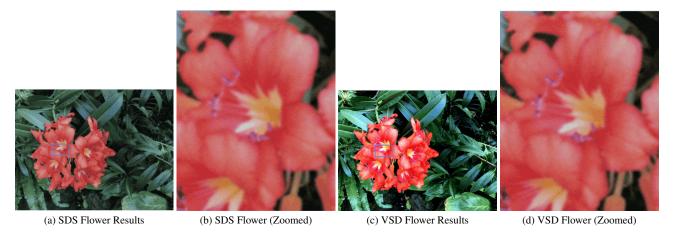


Figure 5. Comparison of SDS and VSD on Flower Images

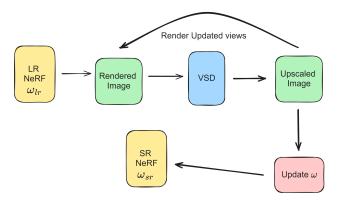


Figure 3. Refer to Algorithm 2

3. Experimental Results

Images of a museum artifact and a flower were generated using RSD, SDS, and VSD. VSD was used to generate images of both the artifact and the flower, whereas RSD was used exclusively for the artifact, and SDS was used solely for the flower. Across all 3 methods, there are noticeable differences in color saturation, texture, and detail.

First, when comparing the museum artifact image gener-

ated with VSD and RSD, there is a clear difference in detail and color saturation. The RSD-generated image has lower color saturation than the VSD-generated image, and furthermore, the VSD-generated image has higher resolution and detail.

Next, comparing the flower images generated with VSD and SDS, there are noticable differences in image distortion and contrast. The SDS-generated image is has less color contrast than the VSD-generated image. Moreover, the VSD image is less distorted and when zoomed in, shows greater detail than the SDS image.

A similar theme exists with figures 12-19. Two VSD-generated images are shown along with zoomed in counterparts, and two RSD-generated photos are shown with zoomed-in counterparts. Comparing VSD 1 (Fig. 14) and RSD 1 (Fig. 16), the VSD-generated image has better resolution and more color saturation. Similarly, when comparing VSD 2 (Fig. 15) with RSD 2 (Fig. 19), there is also more detail and better color saturation in the VSD-generated image than there is in the RSD-generated image.

3.1. Statistical Analysis

We conducted statistical tests to validate the effectiveness of our proposed method, demonstrating significant im-

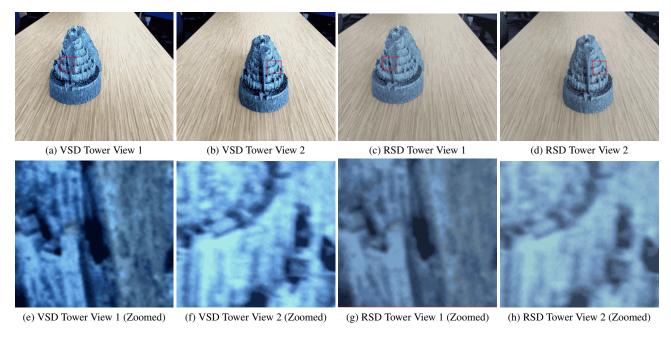


Figure 6. Comparison of VSD and RSD on Tower Images

provements in image resolution, detail clarity, and consistency over existing techniques. The evaluation of the NeRFs generated by our model was based on three standard metrics: LPIPS (Learned Perceptual Image Patch Similarity), NIQE (Natural Image Quality Evaluator), and PSNR (Peak Signal-to-Noise Ratio).

3.2. Comparison with Existing Methods

We compared our method's performance with existing techniques using LPIPS, NIQE, and PSNR scores. Each of these tests had 10000 initial NeRF render steps with Instant-NGP, then 4 rounds of 2000 Super Resolution steps with the mentioned techniques, then we kept 4000 steps of Iterative 3D Synchronization, and arrived at the benchmarks. One interesting finding was that doing spaced training with LoRA, that is training the LoRA every 3 steps instead of 1, it led to slightly better results. The results are summarized in the table below:

Table 1. Comparison of Proposed Method with Existing Methods

Method	LPIPS	NIQE	PSNR
No Changes (Plain RSD)	0.1496	4.983	3.983
With SDS	0.1587	5.667	3.529
With VSD + LoRA Spaced	0.1523	4.457	4.026
With VSD + LoRA	0.1550	4.612	3.998

As we can see, our methods outperform on NIQE and PSNR, showing our images our higher quality and more natural. Where we do lack a little bit is in LPIPS, and this is ex-

plained further in our Limitations section, but we believe it has to do with VSD giving us higher contrast images, which look much different from the ground truth, an issue seen in other projects as well with VSD, such as ProlificDreamer.

3.3. Limitations

While our approach significantly enhances the quality and consistency of Neural Radiance Fields (NeRF) using Variational Score Distillation (VSD) and Low-rank Adaptation (LoRA), it has some limitations.

Specifically, the Learned Perceptual Image Patch Similarity (LPIPS) score, though improved from the baseline, remains lower than that of Renoised Score Distillation (RSD). This indicates that while our model excels in high resolution and view consistency, it may not match RSD in perceptual fidelity. The VSD's emphasis on higher contrast can lead to overemphasis on certain features, potentially compromising perceptual coherence.

In summary, our methodology achieves notable improvements in resolution and consistency but could benefit from further refinement to enhance perceptual quality. Future work should address these limitations to better balance resolution, consistency, and perceptual accuracy.

To add, due to the backpropagation of the LoRA parameters, our method is 15 to 20% slower than RSD, but this is negligible due to advancements in optimization with backpropagation.

4. Conclusion

In this paper, we introduced a new method for improving Neural Radiance Fields (NeRF) with Super Resolution by using Variational Score Distillation (VSD) and I3DS. Our approach significantly enhances the resolution and detail of images while maintaining consistency in the generated outputs.

Our experimental results confirm that our method outperforms traditional techniques. By integrating advanced technologies such as VSD + LoRA and I3DS, we achieved higher quality images that are closer to real-life visuals. The metrics LPIPS, NIQE, and PSNR showed clear improvements in image quality and consistency compared to existing methods.

The practical implications of our findings are vast. They can benefit various applications in 3D modeling, virtual reality, and computer graphics by providing more accurate and realistic images. This advancement could lead to more immersive experiences in gaming and simulations, as well as improved accuracy in professional fields like medical imaging and architectural visualization.

Future research could focus on further refining these techniques and exploring their application in different contexts or with different types of data. Additional work on reducing computational demands and increasing processing speed could make these improvements more accessible for real-time applications.

In conclusion, our method sets a new standard for NeRF editing, promising more realistic and consistent 3D image generation.

References

- [1] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021. 1
- [2] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. arXiv preprint arXiv:2111.12077, 2022. 1
- [3] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision*, pages 333–350. Springer, 2022. 3
- [4] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. arXiv preprint arXiv:2006.11239, 2020. 1, 2
- [5] Robert Horstmeyer, Phillip Isola, and Richard Zhang. Differentiable stereoscopic scene refinement for neural radiance fields. arXiv preprint arXiv:2111.13652, 2021.
- [6] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu Chen. Lora:

- Low-rank adaptation of large language models. *arXiv* preprint arXiv:2106.09685, 2022. 1, 2, 3
- [7] Jie Long Lee, Chen Li, and Gim Hee Lee. Disrnerf: Diffusion-guided view-consistent super-resolution nerf. arXiv preprint arXiv:2404.00874, 2024.
- [8] Andreas Lugmayr, Martin Danelljan, Luc Van Gool, and Radu Timofte. Srflow: Learning the super-resolution space with normalizing flow. arXiv preprint arXiv:2106.01585, 2021. 2, 3
- [9] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF* Conference on Computer Vision and Pattern Recognition, pages 7210–7219, 2021. 1
- [10] 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. arXiv preprint arXiv:2003.08934, 2020. 1
- [11] Anish Mittal, Anush Krishna Moorthy, and Alan Conrad Bovik. Making a "completely blind" image quality analyzer. *IEEE Signal Processing Letters*, 20(3):209–212, 2012. 4
- [12] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics* (*ToG*), 41(4):1–15, 2022. 3
- [13] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. arXiv preprint arXiv:2102.09672, 2021. 1, 2
- [14] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 5865–5874, 2021. 3
- [15] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pages 10318– 10327, 2021. 1
- [16] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 14335– 14345, 2021. 3
- [17] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. arXiv preprint arXiv:2112.10752, 2022. 2
- [18] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. arXiv preprint arXiv:2011.13456, 2020.
- [19] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. arXiv preprint arXiv:2111.11215, 2022. 1
- [20] Matthew Tancik, Vincent Casser, Xinchen Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron,

- and Henrik Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. *arXiv preprint arXiv:2202.05263*, 2022. 1
- [21] Ayush Tewari, Justus Thies, Ben Mildenhall, Pratul Srinivasan, Edgar Tretschk, Yifan Wang, Christoph Lassner, Vincent Sitzmann, Ricardo Martin-Brualla, Stephen Lombardi, et al. Advances in neural rendering. Computer Graphics Forum, 41(2):703–735, 2022.
- [22] Shrey Vishen, Jatin Sarabu, Rithwick Lakshmanan, Chinmay Bharathulwar, and Vishnu Srinivas. Advancing superresolution in neural radiance fields via variational diffusion strategies. *arXiv preprint arXiv:2410.18137*, 2024. 1
- [23] 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. In Advances in Neural Information Processing Systems, volume 34, pages 27171–27183, 2021. 3
- [24] 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 Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 4690–4699, 2021. 4
- [25] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Chen Change Loy, Yu Qiao, and Xiaoou Tang. Esrgan: Enhanced super-resolution generative adversarial networks. arXiv preprint arXiv:1809.00219, 2018. 3
- [26] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation, 2023.
- [27] Yuanbo Xiangli, Linning Xu, Xingang Pan, Nanxuan Zhao, Anyi Rao, Christian Theobalt, Bo Dai, and Dahua Lin. Citynerf: Building nerf at city scale. arXiv preprint arXiv:2112.05504, 2021.
- [28] Yiheng Xie, Towaki Takikawa, Shunsuke Saito, Or Litany, Shiqin Yan, Numair Khan, Federico Tombari, James Tompkin, Vincent Sitzmann, and Srinath Sridhar. Neural fields in visual computing and beyond. In *Computer Graphics Forum*, volume 41, pages 641–676. Wiley Online Library, 2022. 3
- [29] 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. 3
- [30] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2021. 1
- [31] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. Plenoctrees for real-time rendering of neural radiance fields. arXiv preprint arXiv:2103.14024, 2021. 1
- [32] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. In *European Conference on Computer Vision*, pages 405–421. Springer, 2020. 1