# NeRF Inpainting with Geometric Diffusion Prior and Balanced Score Distillation

Menglin Zhang   Xin Luo   Yunwei Lan   Chang Liu   Rui Li   Kaidong Zhang   Ganlin Yang   Dong Liu

University of Science and Technology of China, Hefei, China

{zhangmenglin, xinluo, ywlan}@mail.ustc.edu.cn, dongeliu@ustc.edu.cn

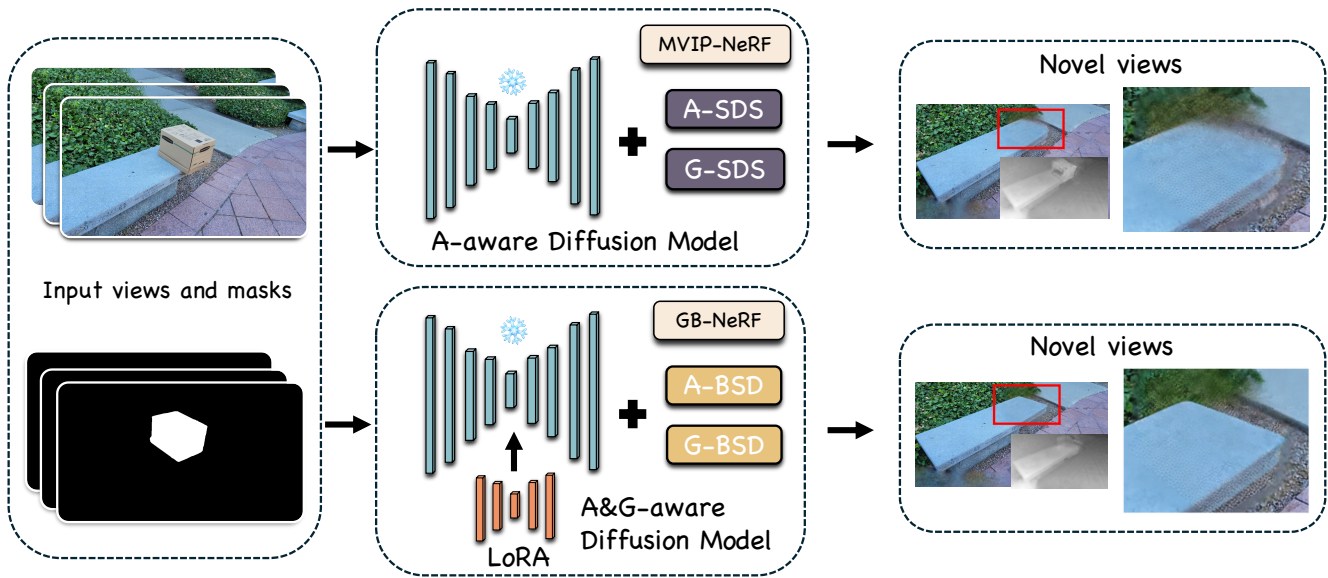**https://github.com/Arcxml/GB-NeRF**

November 26, 2024

Figure 1. Overview of our GB-NeRF framework compared to MVIP-NeRF. Both approaches leverage appearance (A) and geometric (G) priors from diffusion models through score distillation. To enhance geometric accuracy, we introduce two key innovations: (1) a specialized fine-tuning strategy using RGB-normal image pairs; (2) Balanced Score Distillation (BSD), which eliminates high-variability terms present in existing methods like SDS [22] and CSD [37], providing more stable supervision for occluded regions. Compared to MVIP-NeRF [1], our method achieves superior consistency and accuracy in inpainted regions.

## Abstract

*Recent advances in NeRF inpainting have leveraged pretrained diffusion models to enhance performance. However, these methods often yield suboptimal results due to their ineffective utilization of 2D diffusion priors. The limitations manifest in two critical aspects: the inadequate capture of geometric information by pretrained diffusion models and the suboptimal guidance provided by existing Score Distillation Sampling (SDS) methods. To address these problems, we introduce GB-NeRF, a novel framework that enhances NeRF inpainting through improved utilization of 2D diffusion priors. Our approach incorporates two key inno-vations: a fine-tuning strategy that simultaneously learns appearance and geometric priors and a specialized normal distillation loss that integrates these geometric priors into NeRF inpainting. We propose a technique called Balanced Score Distillation (BSD) that surpasses existing methods such as Score Distillation (SDS) and the improved version, Conditional Score Distillation (CSD). BSD offers improved inpainting quality in appearance and geometric aspects. Extensive experiments show that our method provides su-perior appearance fidelity and geometric consistency compared to existing approaches.*

## 1. Introduction

As a pioneering work in neural rendering [7, 17, 20], NeRF (Neural Radiance Field) [17] reconstructs complete 3D scenes from partial viewpoint observations, enabling exceptional 3D reconstruction and novel view synthesis. Practical applications often involve missing areas or unwanted objects that need to be removed. This creates a dual challenge: rendering unobserved viewpoints and filling in missing parts. NeRF inpainting addresses these challenges by reconstructing complete 3D scenes from masked images, ultimately resulting in a full NeRF model of the scene. This technology applies to various areas of 3D content creation, including object removal and scene completion. In this study, we will specifically focus on the task of removing objects.

Traditional NeRF inpainting methods generally follow a two-stage process: first, using 2D inpainting models to complete each view independently, and then employing these inpainted images to train the NeRF model [18]. However, these methods often yield inferior results due to the limited performance of 2D models and a lack of 3D consistency. The emergence of diffusion models [8, 23, 28] has introduced promising solutions. DreamFusion [22] pioneered Score Distillation Sampling (SDS) to leverage 2D diffusion image priors for 3D scene generation. Building on this, MVIP-NeRF employs complementary score distillation losses to enhance both appearance and geometry reconstruction. Despite these advancements, current methods struggle to generate high-quality NeRFs with accurate geometry due to the ineffective use of 2D diffusion priors, which arise from inadequate geometric information capture and suboptimal optimization guidance.

To address these challenges, we introduce GB-NeRF, a novel framework that leverages 2D diffusion priors for high-quality NeRF inpainting. We develop a fine-tuning strategy that trains the model to generate both RGB images and normal maps, utilizing a high-quality RGB-normal image dataset enhanced with captions generated by BLIP. [12]. We also incorporate LoRA [10] into the U-Net [24] and text encoder to learn appearance and geometric information. Additionally, we observe that existing score distillation methods, such as SDS [22] and CSD [37], suffer from unnecessary optimization variability due to random noise and unconditional noise prediction terms. This variability presents a challenge for NeRF inpainting tasks, as consistent supervisory signals are crucial in occluded regions. To mitigate this, we introduce Balanced Score Distillation (BSD), designed specifically for NeRF inpainting tasks, which reduces optimization uncertainty and improves quality in appearance and geometric aspects.

We perform extensive experiments on two representative datasets: *LLFF* [16], a conventional dataset, and *SPIn-NeRF* [18], a more challenging dataset. The results of our experiments demonstrate that our method achieves state-of-the-art performance in both quantitative metrics and visual quality by effectively utilizing diffusion priors. Our primary contributions can be summarized as follows:

1. We propose a specialized diffusion model fine-tuning strategy that improves geometric understanding by learning to generate both RGB images and normal maps together, leading to more effective geometric priors for NeRF inpainting.
2. We introduce Balanced Score Distillation (BSD), a novel optimization technique specifically designed for NeRF inpainting. By removing high-variability terms, BSD offers more stable and consistent supervision signals, enhancing optimization efficiency and improving reconstruction quality.
3. We conduct comprehensive experiments using two representative datasets, demonstrating that our method achieves state-of-the-art performance in both visual quality and geometric accuracy for NeRF inpainting tasks.

## 2. Related Work

### 2.1. NeRF Inpainting

Recent research on NeRF inpainting has primarily focused on addressing 3D inconsistency using restored individual 2D images. Existing studies introduce methods to modify objects represented by NeRFs, such as EditNeRF [14], Clip-NeRF [33], and LaTeRF [19]. A notable advancement is SPIn-NeRF [18], which replaces traditional pixel loss with a more relaxed perceptual loss, enabling NeRF to reveal more high-frequency details while using depth predictions to supervise geometric structure. Another study utilizes LaMa for image inpainting, but its overall quality remains inferior to diffusion models. InpaintNeRF360 [34] employs a similar strategy, relying on perceptual and depth losses to enhance appearance and shape. However, we found that relying solely on perceptual loss does not fundamentally resolve the issues, often leading to subpar results. In contrast, MVIP-NeRF [1] leverages diffusion models as priors to improve appearance and geometric fidelity. We discovered that pretrained diffusion models inadequately capture geometric information, prompting us to fine-tune the model to fully utilize this information.

### 2.2. Learning 3D model via Diffusion Priors

Recently, significant advancements in image generation have been driven by diffusion models [4, 9, 29, 30]. By training on large-scale text-image pairs, these models have achieved remarkable success in text-to-image generation [25], with Stable Diffusion as a notable example. Consequently, many efforts have explored using diffusion models as priors for various image restoration tasks. In addition
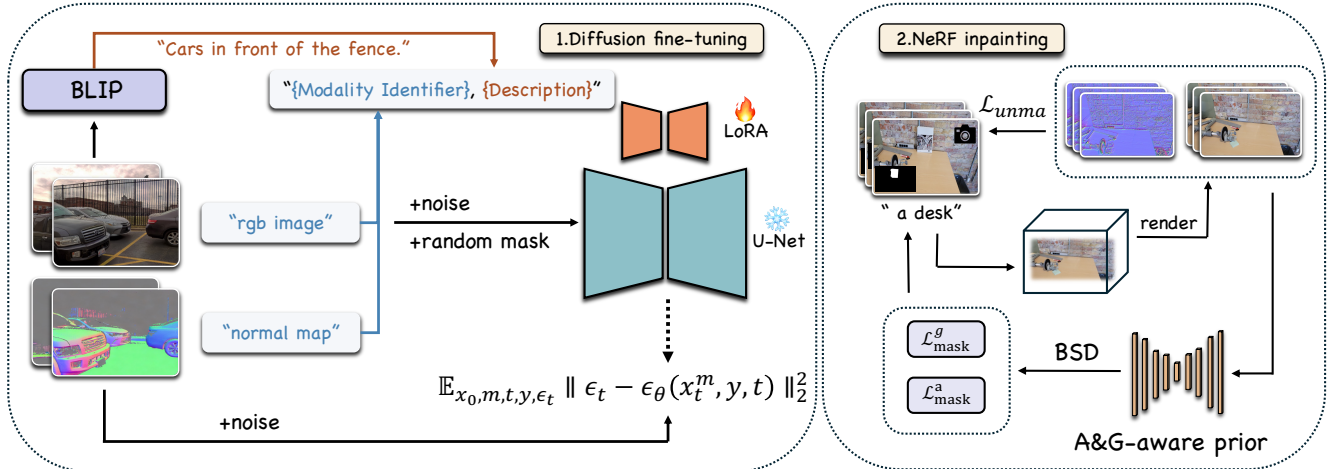
Figure 2. Method overview. (Left) Diffusion fine-tuning: Our approach utilizes the DIODE dataset [32], which provides high-quality RGB images and corresponding normal maps. Captions are generated from RGB images using BLIP [12] and shared with their corresponding normal maps to leverage Stable Diffusion's text understanding capabilities. Modality identifiers ('normal map' or 'RGB image') are prepended to these captions to distinguish between modalities. LoRA is integrated into both U-Net and text encoder to enhance the model's learning capacity further. (Right) NeRF inpainting: Given posed RGB images with corresponding masks and text descriptions, GB-NeRF reconstructs realistic textures and accurate geometry through dual supervision. In unmasked regions, direct pixel-wise RGB reconstruction loss ($L_{unma}$) provides supervision, while in masked areas, our BSD loss guides both RGB image and normal map generation using the fine-tuned diffusion model.

to 2D applications, diffusion priors have gained attention in 3D generation. A pioneering work is Dreamfusion [22], which utilizes multi-view 2D diffusion priors for 3D generation through SDS loss. SDS has been widely adopted in subsequent works [2, 11, 13, 15, 27, 35, 38] aimed at enhancing DreamFusion. For example, Magic3D [13] and Fantasia3D [2] explore optimizing mesh topology for efficient high-resolution rendering. Several methods [6, 26, 39] have also applied SDS to inpainting undesired regions in NeRF scenes, including Nerfiller [36] and MVIP-NeRF [1]. However, these SDS-based methods often struggle to produce high-quality objects. Approaches like CSD [37] and VSD [35], which focus on distillation sampling, still face optimization variability due to random noise and unconditional noise prediction terms. This variability poses challenges for NeRF inpainting tasks, where consistent supervision signals in occluded regions are crucial. To address this, we enhanced the distillation method to provide stable supervision signals in occluded areas, improving generation quality.

## 3. Preliminary

**Neural Radiance Fields.** NeRFs [17] represent a 3D scene using a function $g$ that maps a 3D coordinate $\mathbf{p}$ and a viewing direction $\mathbf{d}$, to a color value $\mathbf{c}$ and a density $\sigma$. Specifically, the function $g$ is a neural network parameterized by $\theta$, where $g_\theta : (\gamma(\mathbf{p}), \gamma(\mathbf{d})) \mapsto (\mathbf{c}, \sigma)$, with $\gamma$ denoting the positional encoding. Each expected pixel color $\hat{C}(\mathbf{r})$

is rendered by casting a ray $\mathbf{r}$ with near and far bounds $t_n$ and $t_f$. The ray segment is typically divided into $N$ intervals $(t_1, t_2, ..., t_N)$, and the pixel color is computed by $\hat{C}(\mathbf{r}) = \sum_{i=1}^{N} w_i \mathbf{c}_i$, where $w_i = T_i(1 - \exp(-\sigma_i \delta_i))$, $T_i = \exp(-\sum_{j=1}^{i-1} \sigma_j \delta_j)$, and $\delta_i = t_i - t_{i-1}$. Mathematically, the NeRF reconstruction loss is formulated as:

$$\mathcal{L}^a = \sum_{\mathbf{r} \in R} ||\hat{C}(\mathbf{r}) - C(\mathbf{r})||^2, \qquad (1)$$

where $\hat{C}(\mathbf{r})$ is the rendered pixel color from the $N$ samples, $R$ is the batch of rays sampled from the training views, and $C(\mathbf{r})$ is the ground truth color for the pixel. When depth information is available, an additional reconstruction loss can be added to optimize the geometry of NeRF scenes [3]:

$$\mathcal{L}^g = \sum_{\mathbf{r} \in R} ||\hat{D}(\mathbf{r}) - D(\mathbf{r})||^2, \qquad (2)$$

where $\hat{D}(\mathbf{r})$ is the rendered depth, and $D(\mathbf{r})$ is the ground-truth depth for the pixel.

**Score Distillation Sampling.** SDS [22] allows the optimization of any differentiable image generator, such as NeRFs or images. Formally, let $\mathbf{x} = g(\theta)$ represent an image rendered by a differentiable generator $g$ with parameter $\theta$. SDS minimizes the density distillation loss [21], which is essential to the KL divergence between the posterior of $\mathbf{x} = g(\theta)$ and the text-conditional density $p_\phi^\omega$:

$$\mathcal{L}_{\text{Dist}}(\theta) = \mathbb{E}_{t,\epsilon}\big[w(t)\, D_{\text{KL}}\big(q(\mathbf{x}_t|\mathbf{x}) \,\|\, p_\phi^\omega(\mathbf{x}_t; y, t)\big)\big], \quad (3)$$

where $w(t)$ is a weighting function, $y$ is the text embedding, and $t$ is the noise level. For efficient computation, SDS updates the parameter $\theta$ by selecting random timesteps $t \sim \mathcal{U}(t_{\min}, t_{\max})$, forwarding $\mathbf{x} = g(\theta)$ with noise $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and calculating the gradient as:

$$\nabla_\theta \mathcal{L}_{\text{SDS}}(\theta) = \mathbb{E}_{t,\epsilon}\left[w(t)\big(\epsilon_\phi^\omega(\mathbf{x}_t; y, t) - \epsilon\big)\frac{\partial \mathbf{x}}{\partial \theta}\right]. \quad (4)$$

## 4. Method

Our task involves a set of RGB images, denoted as $\mathcal{I} = \{I_i\}_{i=1}^n$, along with their corresponding camera poses $\mathcal{G} = \{G_i\}_{i=1}^n$ and object masks $\mathcal{M} = \{m_i\}_{i=1}^n$. The objective is to train a NeRF (Neural Radiance Fields) model to render inpainted scene content from any novel viewpoint. Our approach consists of two main components. First, we introduce our fine-tuning strategy, which allows diffusion models to learn both appearance and geometric priors (see Section § 4.1). Next, we enhance Classifier Score Distillation (CSD) (refer to Section § 4.2.1) and propose Balanced Score Distillation (BSD), an improved alternative to traditional Score Distillation Sampling (SDS) methods (see Section § 4.2.2). Figure 2 illustrates our entire approach.

### 4.1. Geometric Prior Enhancement

Building on the insights from MVIP-NeRF [1], we understand the importance of accurate geometric reconstruction in NeRF inpainting. Although MVIP-NeRF has shown that pretrained Stable Diffusion has inherent geometric priors and can handle normal maps, our experiments indicate that it struggles to effectively complete these normal maps. This challenge arises mainly from diffusion models' tendency to produce overly smooth results, which can compromise the intricate geometric details crucial for accurate 3D reconstruction. To overcome this limitation, we fine-tune the diffusion model to improve its ability to generate structurally accurate normal maps while still being effective in generating RGB images.

The fine-tuning process incorporates Low-Rank Adaptation (LoRA) [10] into both the U-Net and text encoder, preserving Stable Diffusion's powerful image priors and text understanding capabilities. Our training data comes from the DIODE dataset [32], which provides RGB images and their corresponding normal maps. Besides, we use BLIP [12] to generate captions for RGB images, which are then shared with their corresponding normal maps. Each caption starts with a modality identifier—either "normal map" or "RGB image"—to distinguish between them, which is referred to as caption $y$. The training process

follows the methodology outlined in [40], utilizing a self-supervised inpainting loss for both RGB images and normal maps. In this approach, random masks are sampled and combined for each training image.

### 4.2. Balanced Score Distillation (BSD)

Neural Radiance Fields (NeRFs) have significantly advanced the representation of 3D scenes, allowing for view synthesis from various angles. However, NeRFs face challenges in reconstructing masked or occluded regions within multi-view images. To address this issue, Score Distillation Sampling (SDS) and its improved version, Classifier Score Distillation (CSD), utilize pretrained diffusion models to enhance the optimization of NeRFs. This approach has captured our interest.

However, we have observed that CSD struggles to achieve satisfactory results in NeRF inpainting tasks. We hypothesize that effective NeRF inpainting requires more stable supervision signals. To test this hypothesis, we introduce an additional hyperparameter to control the unconditional noise prediction term in CSD's formulation. Our analysis of the relationship between unconditional noise prediction and inpainting quality informs the development of an improved optimization strategy.

### 4.2.1. Classifier Score Distillation (CSD) Refinement

To better understand the limitations of existing methods, we examine the core optimization mechanism in score distillation. Specifically, the gradient applied to the rendered image $\mathbf{x}$ during optimization plays a crucial role. In SDS, this gradient is denoted as $\delta_x(\mathbf{x}_t; y, t) := \epsilon_\phi(\mathbf{x}_t; y, t) - \epsilon$. This formulation encourages the rendered images to concentrate in high-density areas that are conditioned on the text prompt $y$. In practice, however, classifier-free guidance (CFG) is utilized in diffusion models with a large guidance weight $\omega$ (e.g., $\omega = 100$ in DreamFusion [22]) to achieve high-quality results, causing the final gradient applied to the rendered image to deviate from 4. Specifically, with CFG, $\delta_x$ is expressed as:

$$\delta_x(\mathbf{x}_t; y, t) = \underbrace{[\epsilon_\phi(\mathbf{x}_t; y, t) - \epsilon]}_{\delta_x^{\text{gen}}} + \omega \cdot \underbrace{[\epsilon_\phi(\mathbf{x}_t; y, t) - \epsilon_\phi(\mathbf{x}_t; t)]}_{\delta_x^{\text{cls}}}. \quad (5)$$

Previous research [37] has demonstrated that the term $\delta_x^{\text{cls}}$ plays a crucial role in determining the optimization direction. Based on this observation, CSD simplifies the process by concentrating exclusively on this dominant term. By also incorporating negative prompts, CSD utilizes the following gradient:

$$\delta_x^{\text{cls}} = \omega_1 \cdot \epsilon_\phi \left( \mathbf{x}_t; y, t \right) + \left( \omega_2 - \omega_1 \right) \cdot \epsilon_\phi \left( \mathbf{x}_t; t \right) \\ - \omega_2 \cdot \epsilon_\phi \left( \mathbf{x}_t; y_{\text{neg}}, t \right), \quad (6)$$

where the unconditional noise prediction term $\epsilon_\phi \left( \mathbf{x}_t; t \right)$ in the equation is determined by two parameters simultaneously. To enhance its flexibility, we modify the formula as follows:

$$\delta_x^{\text{cls}} = \omega_1 \cdot \epsilon_\phi \left( \mathbf{x}_t; y, t \right) \\ + \omega_3 \cdot \epsilon_\phi \left( \mathbf{x}_t; t \right) - \omega_2 \cdot \epsilon_\phi \left( \mathbf{x}_t; y_{\text{neg}}, t \right), \quad (7)$$

which provides greater flexibility in controlling each term's contribution to overall optimization.

To investigate the impact of unconditional noise prediction on NeRF inpainting, we conduct experiments with varying values of $\omega_3$, as shown in Figure 3. Our results reveal that positive weights lead to blurry reconstructions, while negative weights introduce undesirable artifacts. Notably, the reconstruction quality significantly improves as $\omega_3$ approaches zero. This observation aligns with theoretical expectations, as unconditional predictions inherently exhibit higher diversity compared to conditional ones. These findings suggest that NeRF inpainting benefits from more deterministic score estimation, which motivates our proposed improvement detailed below.

#### 4.2.2. Balanced Score Distillation

Our analysis confirms consistent and stable supervision signals in occluded regions are crucial for high-quality NeRF inpainting. Based on this insight, we propose to simply eliminate the unconditional noise prediction term, resulting in the following formulation:

$$\delta_x^{\text{BSD}} = \omega_1 \cdot \epsilon_\phi \left( \mathbf{x}_t; y, t \right) - \omega_2 \cdot \epsilon\phi \left( \mathbf{x}_t; y_{\text{neg}}, t \right). \quad (8)$$

We refer to this approach as Balanced Score Distillation (BSD) because it ensures a balance between positive and negative prompts. BSD simplifies CSD while achieving superior performance, requiring only two network evaluations instead of three, which reduces computational overhead and simplifies parameter tuning. Compared to SDS, our approach eliminates the random noise term and introduces conditional guidance with negative prompts. This design offers enhanced performance while maintaining algorithmic simplicity: the positive prompt term guides optimization toward desired outcomes, while the negative prompt term steers it away from undesirable results, making BSD particularly effective for NeRF inpainting tasks.

#### 4.2.3. Overall Loss

Specifically, we utilize pixel-wise color (Eq. 1) and depth (Eq. 2) reconstruction loss for unmasked regions, which we denote as $\mathcal{L}_{\text{unma}}^a$ and $\mathcal{L}_{\text{unma}}^g$. For masked regions, we first
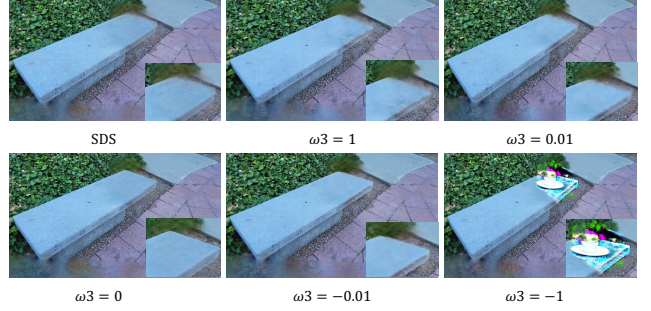


Figure 3. Impact of tuning coefficient $\omega_3$ on NeRF inpainting. The incorporation of unconditional noise prediction introduces excessive randomness, resulting in degraded inpainting quality.

render an RGB image and a normal map from the NeRF scene. We employ our BSD losses to compute a gradient direction iteratively for detailed and high-quality appearance and geometry completion within the mask. The loss in the mask of appearance can be written as:

$$\nabla_\theta \mathcal{L}_{\text{masked}}^a = \\ \left[ \omega_1 \cdot \epsilon_\phi \left( \mathbf{x}_t; y, t \right) - \omega_2 \cdot \epsilon_\phi \left( \mathbf{x}_t; y_{\text{neg}}, t \right) \right] \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial \theta}. \quad (9)$$

Similarly, the loss in the mask of geometry can be written as:

$$\nabla_\theta \mathcal{L}_{\text{masked}}^g = \\ \left[ \omega_1 \cdot \epsilon_\phi \left( \mathbf{n}_t; y, t \right) - \omega_2 \cdot \epsilon_\phi \left( \mathbf{n}_t; y_{\text{neg}}, t \right) \right] \frac{\partial \mathbf{n}_t}{\partial \mathbf{n}} \frac{\partial \mathbf{n}}{\partial \theta}. \quad (10)$$

In general, we address unmasked and masked regions separately, following this general formulation:

$$\mathcal{L} = \mathcal{L}_{\text{unma}}^a + \lambda_1 \mathcal{L}_{\text{unma}}^g + \lambda_2 \mathcal{L}_{\text{mask}}^a + \lambda_3 \mathcal{L}_{\text{mask}}^g. \quad (11)$$

## 5. Experiments and Discussions

**Implementation Details.** We implement our NeRF inpainting model based on SPIn-NeRF [18] and train it on a single NVIDIA A100 GPU for 10,000 iterations using the Adam optimizer with a learning rate of $10^{-4}$. For the diffusion prior, we set the size of all latent inputs to $256 \times 256$ and configure the timestep range with $t_{\min} = 0.02$ and $t_{\max} = 0.98$. For the loss of masked region, we set $\omega_1 = 7.5$ and $\omega_2 = 6.5$ for appearance BSD $\mathcal{L}_{\text{masked}}^a$, while using $\omega_1 = 1.5$ and $\omega_2 = 0.5$ for geometric BSD $\mathcal{L}_{\text{masked}}^g$. For the balance weights in the final loss function (Eq. 11), we empirically set $\lambda_1 = 0.1$ and $\lambda_2 = \lambda_3 = 0.0001$. For diffusion fine-tuning, we set the learning rate to $10^{-4}$ and use LoRA of rank 32.
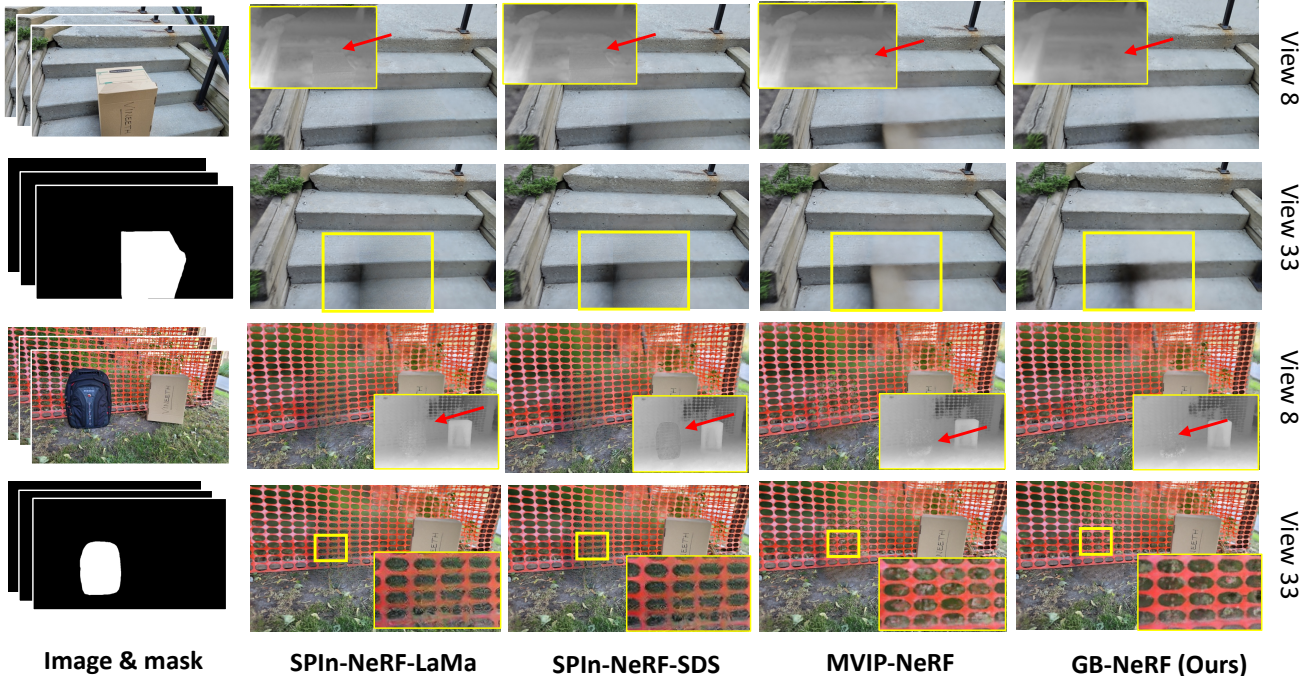
Figure 4. Visual comparison with three representative approaches on two scenes. The first scene uses the prompt 'A stair' while the second uses 'A fence'. Our method effectively handles both scenarios, producing view-consistent results with superior geometric accuracy (note the well-preserved stair structure and depth, while (a) and (b) preserve residual geometry from original objects, and (c) generates artifacts) and realistic textures (observe the cleaner and more structured fence pattern in our results).

Table 1. **Quantitative comparison on *SPIN-NeRF* and *LLFF* datasets.** Our method demonstrates superior performance across most metrics compared to existing NeRF inpainting approaches.

| | SPIN-NeRF | | | | | | | LLFF | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | FID↓ | NIMA↑ | BRISQUE↓ | D-FID↓ | D-PSNR↑ | FID↓ | NIMA↑ | BRISQUE↓ |
| SPIn-NeRF + LaMa | 19.651 | 0.4197 | 79.424 | 4.081 | 22.433 | 186.967 | 13.827 | 283.348 | 4.788 | 12.325 |
| SPIn-NeRF + SDS | 19.655 | 0.4199 | 79.196 | 4.101 | 22.420 | 191.496 | 13.834 | 286.742 | 4.729 | 12.345 |
| MVIP-NeRF | **19.813** | 0.4208 | 72.616 | 4.455 | 23.562 | 172.127 | 13.914 | 285.873 | 4.862 | 11.901 |
| GB-NeRF (Ours) | 19.489 | **0.4266** | **67.587** | **4.550** | **17.496** | **150.473** | **14.157** | **270.923** | **4.915** | **11.711** |

Table 2. **Ablation analysis.** Quantitative comparison of different model components. Enhanced geometric priors significantly improve geometric-related metrics (D-FID), while Balanced Score Distillation (BSD) enhances overall visual quality by reducing optimization uncertainty. The best and second-best results are highlighted in **bold** and underlined, respectively.

| | PSNR↑ | SSIM↑ | FID↓ | NIMA↑ | BRISQUE↓ | D-FID↓ |
|---|---|---|---|---|---|---|
| origin | **19.814** | <u>0.421</u> | 72.616 | 4.455 | 23.562 | 172.127 |
| +LoRA | 19.543 | 0.402 | 70.137 | 4.432 | **17.029** | **147.330** |
| +BSD | <u>19.813</u> | 0.419 | **64.310** | <u>4.550</u> | 23.002 | 161.416 |
| +LoRA+BSD | 19.489 | **0.426** | <u>67.587</u> | **4.551** | <u>17.496</u> | <u>150.473</u> |

**Datasets.** We evaluate our method on two real-world datasets: *SPIN-NeRF* and *LLFF*. Below, we provide detailed descriptions of these datasets.

*SPIn-NeRF* [18] serves as an object removal benchmark consisting of 10 scenes. For each scene, the dataset provides 60 training views captured with an object intended for removal, accompanied by corresponding inpainting masks indicating the object's location. For evaluation purposes, we select 8 challenging scenes from this dataset, where each scene includes 40 testing views in which the target object has been physically removed during capture.

*LLFF* [16] comprises multiple real-world scenes with varying numbers of images (ranging from 20 to 45). For our experiments, we utilize a four-scene subset that has been annotated with 3D grounded object removal masks. Since this dataset does not provide a separate test set, we follow previous methods and use the training views for evaluation purposes. Following established practices in prior work [1],
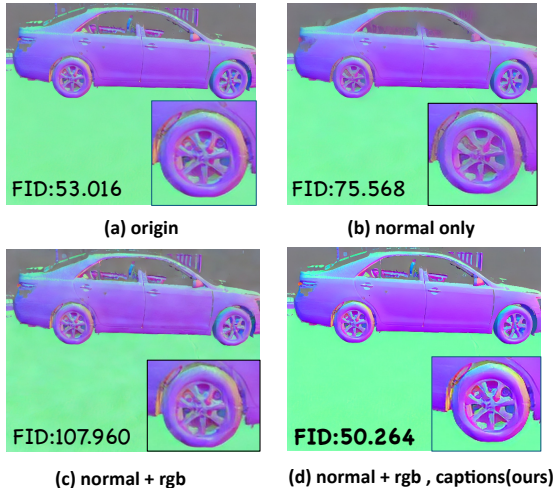
Figure 5. Comparison of different fine-tuning strategies for the diffusion model. (a) Original diffusion model without fine-tuning; (b) Fine-tuning with normal maps only, using modality identifier 'normal map' as prompt; (c) Fine-tuning with both RGB images and normal maps, using modality identifiers 'RGB image' and 'normal map' as prompts; (d) Our approach: fine-tuning with RGB-normal image pairs using BLIP-generated captions prepended with modality identifiers as prompts. Results show that while strategies (b) and (c) underperform compared to the original model, our method significantly enhances the model's capability in normal map reconstruction.

we standardize all images by resizing them to have a long-edge size of 1008.

**Metrics.** To comprehensively evaluate our inpainting results, we use several complementary metrics. For direct comparison with ground-truth scenes, we utilize traditional full-reference metrics PSNR and SSIM, which measure pixel-level and structural similarities, respectively. To assess the perceptual quality of generated images, we employ no-reference metrics NIMA and BRISQUE, which evaluate image aesthetics and quality without requiring ground truth. Additionally, we use FID to measure the distribution distance between inpainting results and ground-truth scenes. For evaluating geometric accuracy, we employ D-PSNR to compare the reconstructed disparity maps with the ground-truth disparity maps. Additionally, we report D-FID to measure geometric consistency by calculating FID on the disparity maps.

## 5.1. Results

**Baselines.** Recent methods based on NeRF for inpainting have shown better performance than traditional image and video inpainting techniques. Therefore, we focus our comparison on state-of-the-art NeRF inpainting methods. Specifically, we compare against SPIn-NeRF [18], SPIn-NeRF-SDS, and MVIP-NeRF [1]. SPIn-NeRF applies the

2D inpainting model LaMa [31] to independently inpaint each image in the dataset before training NeRF with these inpainted images. SPIn-NeRF-SDS enhances the NeRF optimization process by applying additional supervision in masked regions using SDS. MVIP-NeRF improves geometric reconstruction by utilizing both an appearance SDS and a geometric SDS, which take RGB images and normal maps as input, respectively. To ensure a fair comparison, we use the official implementations provided by the authors and report the results according to standard evaluation protocols.

**Qualitative analysis.** We present visual comparisons of different methods on the SPIn-NeRF dataset, as shown in Figure 4. The first scene showcases the complex geometric structures of a staircase, which challenge the methods' ability to reconstruct its intricate details accurately. The results indicate that both SPIn-NeRF-based methods produce continuous but blurry areas that do not blend seamlessly with the surrounding scene. This limitation primarily stems from their heavy reliance on LaMa's inpainting results, which lack sufficient quality for complex geometric structures. Despite incorporating SDS priors, SPIn-NeRF-SDS cannot fully overcome the limitations inherited from LaMa's initial processing. While MVIP-NeRF generates smoother transitions, it struggles to maintain sharp geometric edges. In contrast, our method successfully reconstructs highly realistic geometric structures while preserving fine details. The second scene presents a different challenge with an orange net, requiring accurate reproduction of its periodic texture. Among all methods, only MVIP-NeRF and our approach successfully reconstruct the net's repeated pattern. However, our method demonstrates superior performance by producing clearer and more realistic geometric details compared to MVIP-NeRF's results.

**Quantitative evaluation.** As shown in Table 1, our method achieves superior performance across most evaluation metrics. Specifically, our approach demonstrates excellent structural similarity and geometric faithfulness, as evidenced by high SSIM scores and favorable D-PSNR and D-FID metrics. Additionally, our method achieves outstanding visual quality, as indicated by strong performance in FID, NIMA, and BRISQUE metrics. While our method shows relatively lower performance in terms of PSNR, we argue that this metric alone is not sufficiently reliable for evaluating inpainting quality. This limitation of PSNR is particularly relevant given that inpainting is an inherently ill-posed problem [5], where multiple plausible solutions exist for unobserved regions. In such cases, the posterior mean solution usually produces the highest PSNR scores, but it often results in undesirably blurry outputs that compromise visual quality.

**Ablation Studies.** We start by performing ablation studies to validate our key contributions: the enhanced geometric priors in diffusion models and the improved score
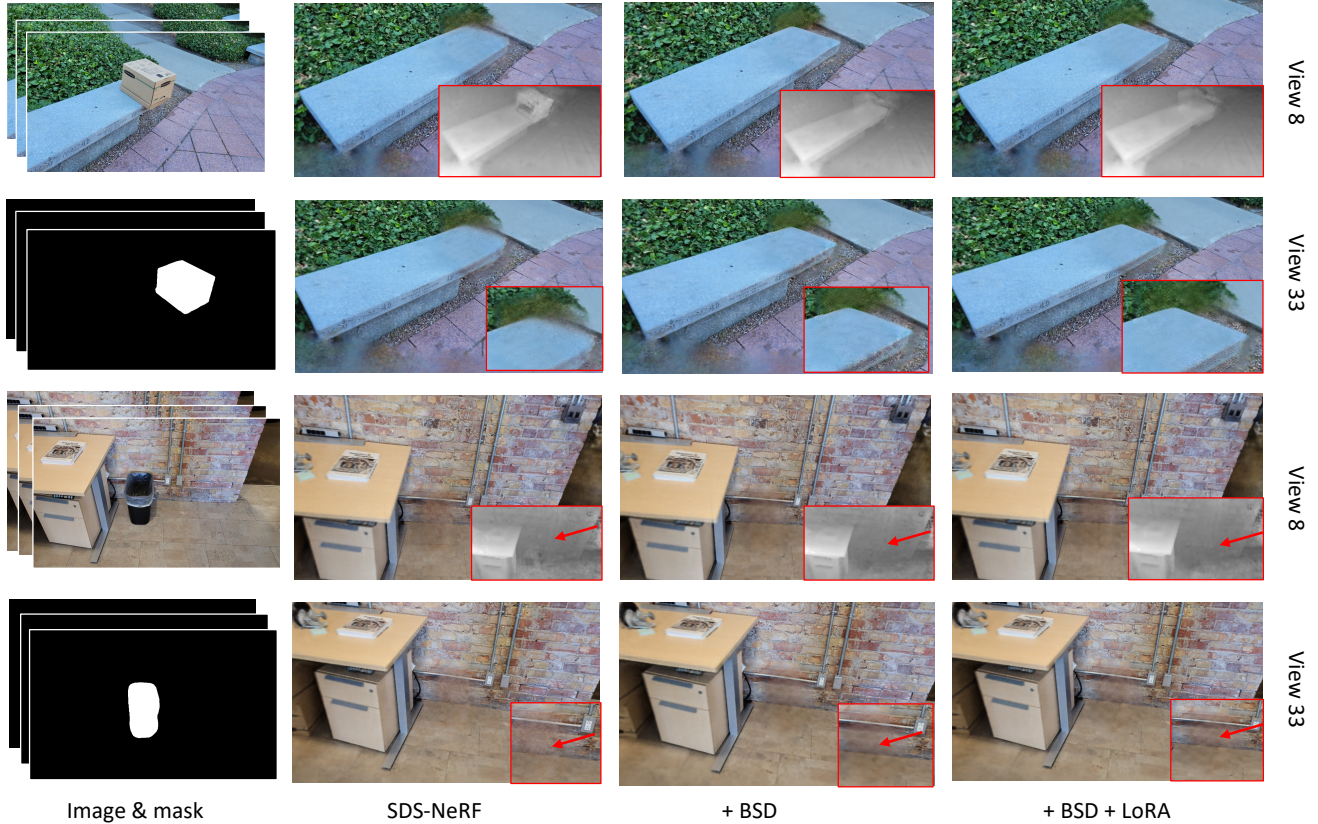
Figure 6. Visual comparison of ablation study about each component. After adding BSD distillation, the blurriness around the edges is reduced (see the bench and the floor gaps), and further incorporating LoRA better reconstructs the geometry of the bench and the wall (see the depth map of the bench and floor).

distillation method. As demonstrated in Table 2, our complete method outperforms others across most metrics. In particular, by integrating LoRA for geometric prior learning, our approach significantly enhances the model's capability for geometric reconstruction, reducing D-FID by 25. Meanwhile, BSD consistently enhances visual quality by reducing optimization uncertainty, as evidenced by FID, NIMA, BRISQUE, and D-FID metrics improvements. Additionally, as discussed before, we ignore PSNR because it is not a reliable metric. By combining these two approaches, we can take advantage of their complementary benefits, resulting in the best overall performance. These quantitative improvements are also reflected in the visual results shown in Fig. 6, where LoRA notably improves geometric reconstruction accuracy, while BSD effectively reduces edge blurriness in the final results. We further validate our strategy for fine-tuning the diffusion model. Our approach is based on two key insights: (1) it is essential to preserve the original appearance prior while also acquiring new geometric priors, and (2) utilizing Stable Diffusion's text comprehension capabilities is advantageous. To verify these insights, we conduct three comparative experiments: using only normal maps with the modality identifier "normal map" as prompt; using both normal maps and RGB images with simple modality identifiers ("RGB image" or "normal map") as prompts; and our approach of using both normal maps and RGB images with descriptive captions prepended with modality identifiers as prompts.

As shown in Fig. 5, alternative approaches yield suboptimal results, highlighting the importance of incorporating both RGB images and descriptive captions in the fine-tuning process.

## 6. Conclusion

In this work, we introduce GB-NeRF, a novel framework that enhances NeRF inpainting by more effectively utilizing 2D diffusion priors. First, we fine-tune the diffusion model to improve its ability to generate structurally accurate normal maps while still producing effective RGB images. Additionally, our proposed Balanced Score Distillation (BSD) technique outperforms existing methods, such as SDS and CSD, by delivering higher-quality inpainting with fewer ar-

tifacts. Our experiments demonstrate that GB-NeRF excels in both appearance fidelity and geometric consistency. However, our work has several limitations: (i) using fine-tuned diffusion priors for geometric enhancement increases training time; (ii) our method introduces new hyperparameters that require efforts to adjust; and (iii) similar to previous research [18], our method cannot eliminate shadows.

# References

[1] Honghua Chen, Chen Change Loy, and Xingang Pan. Mvip-nerf: Multi-view 3d inpainting on nerf scenes via diffusion prior. In *CVPR*, 2024. 1, 2, 3, 4, 6, 7

[2] Rui Chen, Yongwei Chen, Ningxin Jiao, and Kui Jia. Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. *arXiv preprint arXiv:2303.13873*, 2023. 3

[3] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12882–12891, 2022. 3

[4] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *NeurIPS*, volume 34, pages 8780–8794, 2021. 2

[5] Omar Elharrouss, Noor Almaadeed, Somaya Al-Maadeed, and Younes Akbari. Image inpainting: A review. *Neural Processing Letters*, 51:2007–2028, 2020. 7

[6] Ayaan Haque, Matthew Tancik, Alexei A Efros, Aleksander Holynski, and Angjoo Kanazawa. Instruct-nerf2nerf: Editing 3d scenes with instructions. *arXiv preprint arXiv:2303.12789*, 2023. 3

[7] Philipp Henzler, Niloy J Mitra, and Tobias Ritschel. Learning a neural 3d texture space from 2d exemplars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8356–8364, 2020. 2

[8] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 2

[9] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 2

[10] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. 2, 4

[11] Yukun Huang, Jianan Wang, Yukai Shi, Xianbiao Qi, Zheng-Jun Zha, and Lei Zhang. Dreamtime: An improved optimization strategy for text-to-3d content creation. *arXiv preprint arXiv:2306.12422*, 2023. 3

[12] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*, 2022. 2, 3, 4

[13] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,

[14] Steven Liu, Xiuming Zhang, Zhoutong Zhang, Richard Zhang, Jun-Yan Zhu, and Bryan Russell. Editing conditional radiance fields. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5773–5783, 2021. 2

[15] Gal Metzer, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. Latent-nerf for shape-guided generation of 3d shapes and textures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12663–12673, 2023. 3

[16] 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

[17] 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*, volume 12346, pages 405–421. Springer, 2020. 2, 3

[18] Ashkan Mirzaei, Tristan Aumentado-Armstrong, Konstantinos G Derpanis, Jonathan Kelly, Marcus A Brubaker, Igor Gilitschenski, and Alex Levinshtein. Spin-nerf: Multiview segmentation and perceptual inpainting with neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20669–20679, 2023. 2, 5, 6, 7, 9

[19] Ashkan Mirzaei, Yash Kant, Jonathan Kelly, and Igor Gilitschenski. Laterf: Label and text driven object radiance fields. In *European Conference on Computer Vision*, pages 20–36. Springer, 2022. 2

[20] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3504–3515, 2020. 2

[21] Aaron Oord, Yazhe Li, Igor Babuschkin, Karen Simonyan, Oriol Vinyals, Koray Kavukcuoglu, George Driessche, Edward Lockhart, Luis Cobo, Florian Stimberg, et al. Parallel wavenet: Fast high-fidelity speech synthesis. In *International conference on machine learning*, pages 3918–3926. PMLR, 2018. 3

[22] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 1, 2, 3, 4

[23] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2

[24] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015. 2

[25] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo

Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. Laion-5b: An open large-scale dataset for training next generation image-text models. *ArXiv*, abs/2210.08402, 2022. 2

[26] Ruizhi Shao, Jingxiang Sun, Cheng Peng, Zerong Zheng, Boyao Zhou, Hongwen Zhang, and Yebin Liu. Control4d: Dynamic portrait editing by learning 4d gan from 2d diffusion-based editor. *arXiv preprint arXiv:2305.20082*, 2023. 3

[27] Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. Mvdream: Multi-view diffusion for 3d generation, 2023. 3

[28] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015. 2

[29] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, page 2256–2265, 2015. 2

[30] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021. 2

[31] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky. Resolution-robust large mask inpainting with fourier convolutions. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2149–2159, 2022. 7

[32] Igor Vasiljevic, Nick Kolkin, Shanyi Zhang, Ruotian Luo, Haochen Wang, Falcon Z. Dai, Andrea F. Daniele, Mohammadreza Mostajabi, Steven Basart, Matthew R. Walter, and Gregory Shakhnarovich. DIODE: A Dense Indoor and Outdoor DEpth Dataset. *CoRR*, abs/1908.00463, 2019. 3, 4

[33] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Clip-nerf: Text-and-image driven manipulation of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3835–3844, 2022. 2

[34] Dongqing Wang, Tong Zhang, Alaa Abboud, and Sabine Süsstrunk. Inpaintnerf360: Text-guided 3d inpainting on unbounded neural radiance fields. *arXiv preprint arXiv:2305.15094*, 2023. 2

[35] 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. *arXiv preprint arXiv:2305.16213*, 2023. 3

[36] Ethan Weber, Aleksander Holynski, Varun Jampani, Saurabh Saxena, Noah Snavely, Abhishek Kar, and Angjoo Kanazawa. Nerfiller: Completing scenes via generative 3d inpainting. In *CVPR*, 2024. 3

[37] Xin Yu, Yuan-Chen Guo, Yangguang Li, Ding Liang, Song-Hai Zhang, and Xiaojuan Qi. Text-to-3d with classifier score distillation. *arXiv preprint arXiv:2310.19415*, 2023. 1, 2, 3, 4

[38] Joseph Zhu and Peiye Zhuang. Hifa: High-fidelity text-to-3d with advanced diffusion guidance. *arXiv preprint arXiv:2305.18766*, 2023. 3

[39] Jingyu Zhuang, Chen Wang, Lingjie Liu, Liang Lin, and Guanbin Li. Dreameditor: Text-driven 3d scene editing with neural fields. *arXiv preprint arXiv:2306.13455*, 2023. 3

[40] Junhao Zhuang, Yanhong Zeng, Wenran Liu, Chun Yuan, and Kai Chen. A task is worth one word: Learning with task prompts for high-quality versatile image inpainting, 2023. 4