# GaussianSR: 3D Gaussian Super-Resolution with 2D Diffusion Priors

**Xiqian Yu** [*][1] , **Hanxin Zhu** [*][1], **Tianyu He** [2], **Zhibo Chen** [1]
[1] University of Science and Technology of China, [2] Microsoft Research Asia
{yuxiqian,hanxinzhu}@mail.ustc.edu.cn
tianyuhe@microsoft.com
chenzhibo@ustc.edu.cn

## Abstract

Achieving high-resolution novel view synthesis (HRNVS) from low-resolution input views is a challenging task due to the lack of high-resolution data. Previous methods optimize high-resolution Neural Radiance Field (NeRF) from low-resolution input views but suffer from slow rendering speed. In this work, we base our method on 3D Gaussian Splatting (3DGS) due to its capability of producing high-quality images at a faster rendering speed. To alleviate the shortage of data for higher-resolution synthesis, we propose to leverage off-the-shelf 2D diffusion priors by distilling the 2D knowledge into 3D with Score Distillation Sampling (SDS). Nevertheless, applying SDS directly to Gaussian-based 3D super-resolution leads to undesirable and redundant 3D Gaussian primitives, due to the randomness brought by generative priors. To mitigate this issue, we introduce two simple yet effective techniques to reduce stochastic disturbances introduced by SDS. Specifically, we 1) shrink the range of diffusion timestep in SDS with an annealing strategy; 2) randomly discard redundant Gaussian primitives during densification. Extensive experiments have demonstrated that our proposed GaussainSR can attain high-quality results for HRNVS with only low-resolution inputs on both synthetic and real-world datasets. Project page: https://chchnii.github.io/GaussianSR/.

## 1 Introduction

Novel View Synthesis (NVS) has been extensively studied in computer vision and graphics. In particular, Neural Radiance Field (NeRF) [1] has demonstrated its impressive ability to generate high-quality visual content. More recently, 3D Gaussian Splatting (3DGS) [2] has been attracting widespread attention due to its capability of producing high-quality images with faster rendering speed. However, achieving high-resolution novel view synthesis (HRNVS) from low-resolution inputs remains an under-explored yet challenging task.

There exist two primary difficulties for HRNVS. Firstly, previous works [3, 4, 5] mainly rely on optimizing high-resolution NeRF from low-resolution input views. Although these methods can synthesize satisfactory high-resolution novel views, the stratified sampling required for rendering in NeRF is costly and results in high rendering time. Secondly, we only have low-resolution input views to produce high-resolution results. To tackle this, NeRF-SR [3] exploits a supersampling strategy to estimate color and density at the sub-pixel level. However, it is still a challenge to get enough information from the low-resolution input alone.

In this work, we propose GaussianSR, which aims to introduce 2D generative priors learned from large-scale image data into HRNVS. Specifically, we build our method upon 3DGS due to its photo-realistic visual quality and real-time rendering. In order to leverage 2D priors, we derive inspiration

---

from DreamFusion [6], a method that distills 2D diffusion priors into text-to-3D generation with Score Distillation Sampling (SDS). In this way, to introduce 2D priors to HRNVS, a straightforward solution is to distill off-the-shelf 2D super-resolution diffusion priors into 3DGS for high-resolution novel view synthesis. Nevertheless, we notice that applying SDS directly fails with some undesirable and redundant 3D Gaussian primitives. We suspect that this is due to the inherent randomness of the generative priors, as it always takes random noise and timestep as input to produce natural image distribution. This property is particularly amplified in SDS when we aim to optimize high-resolution 3DGS with denser Gaussian primitives, since there are large variances in the gradients during 3DGS densification (a process that clones or splits current Gaussian primitives into more). To mitigate this issue, we propose two simple yet effective techniques, which reduce stochastic disturbances introduced by SDS. Firstly, to alleviate the randomness of the diffusion timestep, we shrink the sampling range of the diffusion timestep with an annealing strategy. Secondly, to prevent explosive Gaussian primitives, we randomly discard redundant primitives during the process of densification. We validate our GaussianSR in various scenarios, including synthesized and realistic scenarios, and experimental results demonstrate that the rendering quality of GaussianSR outperforms existing state-of-the-art methods.

In conclusion, our contributions can be summarized as follows:

- To alleviate the lack of high-resolution data, we, for the first time, propose to distill generative priors of 2D super-resolution models into HRNVS.

- We observe that applying SDS directly to Gaussian-based 3D super-resolution leads to undesirable and redundant 3D Gaussian primitives, due to randomness brought by generative priors. To solve this issue, we propose two techniques to reduce stochastic disturbances introduced by SDS.

- Experimental results demonstrate that our proposed GaussianSR achieves higher-quality HRNVS than the state-of-the-art solutions from only low-resolution inputs.

## 2    Related Work

### 2.1    3D Gaussian Splatting

3D Gaussian Splatting (3DGS) [2] provides a promising and effective approach for Novel View Synthesis (NVS). By representing the scene explicitly as a collection of 3D Gaussian primitives and rendering views through rasterization, 3DGS achieves impressive quality and speed in NVS. This has led to the emergence of various 3DGS extensions, including modeling dynamic scenes [7, 8, 9, 10], reconstruction without input camera poses [11, 12], few-shot view synthesis [13, 14, 15] and applications in other fields [16, 17, 18, 19, 20]. Additionally, some studies [21, 22, 23, 24] focus on advancing 3DGS itself to improve the quality of NVS. However, existing relevant works have primarily concentrated on synthesizing novel views with resolutions limited to the input views, neglecting the high-resolution novel views synthesis (HRNVS) from low-resolution inputs. In this paper, we explore high-quality HRNVS with 3DGS.

### 2.2    High-Resolution Novel View Synthesis

High-resolution novel view synthesis (HRNVS) aims to synthesize high-resolution novel views from only low-resolution inputs. As a pioneer, NeRF-SR [3] optimizes high-resolution NeRF with the sub-pixel constraint, ensuring that the values of low-resolution (LR) pixels equal the mean value of high-resolution (HR) sub-pixels RefSR-NeRF [5] reconstructs the high-frequency details with the help of a single high-resolution reference image. Furthermore, Super-NeRF [4] constructs a consistency-controlling super-resolution module to generate view-consistent high-resolution details for NeRF. However, these NeRF-based methods suffer from slow rendering speed. Recently, 3DGS has gained popularity due to its primitive-based representation, which can produce high-quality images at faster rendering speeds. While a concurrent work, SRGS [25], similarly focuses on 3DGS-based HRNVS, however, our study and SRGS differ significantly not only in terms of technical contributions but also in motivation: we aim to introduce 2D diffusion priors, which is learned from large-scale 2D data, into HRNVS.
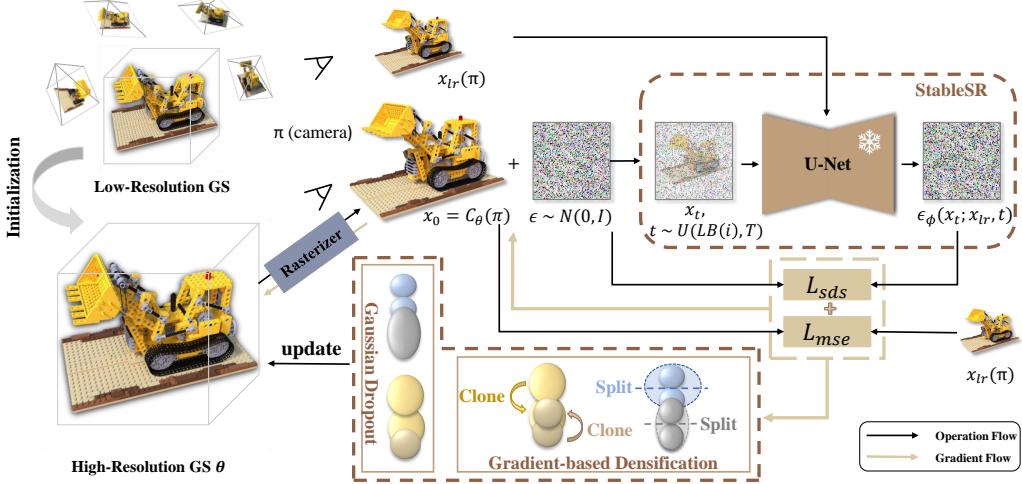
Figure 1: Overview of GaussianSR. To alleviate the lack of high-resolution data, we synthesize high-resolution novel views by distilling 2D diffusion priors into 3D representation with SDS (Sec. 3.1). Since the redundant Gaussian primitives are introduced due to the randomness of generative priors (Sec. 3.2), we propose Gaussian Dropout and diffusion timestep annealing to reduce stochastic disturbance (Sec. 3.3).

## 3 Methodology

In this section, we provide a comprehensive overview of our proposed GaussianSR. To begin with, recognizing the challenge posed by the limit availability of high-resolution data, we leverage 2D diffusion priors distilled by SDS [6] to optimize high-resolution 3DGS (Sec. 3.1). However, the randomness introduced by generative priors will lead to undesirable and redundant Gaussian primitives (Sec. 3.2). Hence, we propose two simple yet effective techniques to mitigate this issue (Sec. 3.3).
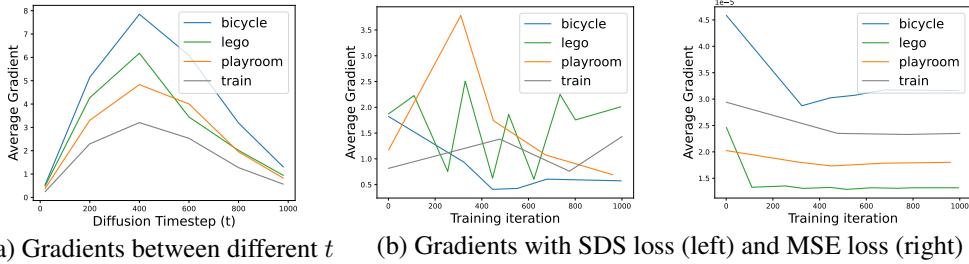
### 3.1 3DGS Super-Resolution with SDS Optimization

**3DGS.** As an effective method for novel view synthesis, 3D Gaussian Splatting (3DGS) [2] represents a 3D scene using a series of Gaussian primitives comprised of position $\boldsymbol{\mu} \in \mathbb{R}^3$, scaling $\boldsymbol{s} \in \mathbb{R}^3$, rotation $\boldsymbol{r} \in \mathbb{R}^3$, and color $\boldsymbol{c} \in \mathbb{R}^3$. To faithfully reconstruct the 3D scene, these Gaussian primitives are initialized with sparse point clouds estimated by SfM [26], followed by a densification operation, *i.e.*, the split and clone operation, that adaptively control their numbers and densities. Concretely, whether a 3D Gaussian primitive is split or cloned is determined by the average gradient magnitude of the Normalized Device Coordinates (NDC) [27] for the viewpoints in which the Gaussian primitive participates in the calculation. For example, for Gaussian primitive $k$ under viewpoint $M_i$, the NDC is $(\mu_{ndc,x}^{k,M_i}, \mu_{ndc,y}^{k,M_i}, \mu_{ndc,z}^{k,M_i})$, and the loss under viewpoint $M_i$ is $\mathcal{L}_{M_i}$. During optimization, Gaussian primitive $k$ participates in the calculation for $M$ viewpoints. When Gaussian primitive satisfies:

$$|g| = \frac{\sum_{M_i=1}^{M} \sqrt{\left(\frac{\partial \mathcal{L}_{M_i}}{\partial \mu_{ndc,x}^{k,M_i}}\right)^2 + \left(\frac{\partial \mathcal{L}_{M_i}}{\partial \mu_{ndc,y}^{k,M_i}}\right)^2}}{M} > \tau_{pos}, \tag{1}$$

it is transformed into two Gaussian primitives, where $\tau_{pos}$ is the default threshold.

**Distilling 2D Diffusion Priors for 3DGS Super-Resolution.** The primary challenge encountered in high-resolution novel view synthesis (HRNVS) from low-resolution inputs is the scarcity of data, a limitation pervasive across various domains. For instance, in text-to-3D generation, the performance of early endeavours [28, 29, 30] is limited by the small-scale text-3D datasets adopted (*e.g.*, ShapeNet [31]), resulting in poor generalization. To overcome the bottleneck of data scarcity

3

(a) Gradients between different $t$     (b) Gradients with SDS loss (left) and MSE loss (right)

Figure 2: **(a)** The gradient values under the constraint of SDS loss are visualized, revealing substantial variance across different diffusion timesteps $t$. **(b)** When comparing the gradient values under two different constraints—SDS loss on the left and MSE loss on the right—the gradient variance for SDS is significantly larger than that for MSE.

and facilitate the generation of more diverse 3D assets, DreamFusion [6] introduces Score Distillation Sampling (SDS), which aims to distill generative priors from pretrained text-to-image diffusion models. Drawing inspiration from DreamFusion, we propose leveraging off-the-shelf 2D super-resolution diffusion priors to mitigate the data shortage challenge in HRNVS.

Specifically, as shown in Fig. 1, given a set of multi-view low-resolution images $x_{lr}$, our objective is to synthesize high-resolution novel views through optimizing high-resolution 3DGS with SDS. Initially, we reconstruct a low-resolution 3DGS from the multi-view low-resolution inputs, which serves as the initialization for the high-resolution 3DGS. Subsequently, we optimize the high-resolution 3DGS using priors distilled from a diffusion-based 2D super-resolution model along with the low-resolution inputs. Let $C_{\boldsymbol{\theta}}(\pi)$ represent the rendered high-resolution image at the given viewpoint $\pi$, where $C$ is the differentiable rendering function for the high-resolution 3DGS parameterized by $\boldsymbol{\theta}$. Our goal is to optimize the rendered high-resolution image, denoted as $x_0 := C_{\theta}(\pi)$, by introducing the SDS loss $\mathcal{L}_{SDS}$, which encourages $x_0$ move toward higher density region conditioned on its corresponding low-resolution image $x_{lr}$. Particularly, $\mathcal{L}_{SDS}$ computes the difference of predicted noise $\epsilon_{\phi}$ and the added noise $\epsilon$ as per-pixel gradient, which is then used to update the high-resolution 3DGS parameters $\theta$:

$$\nabla_{\boldsymbol{\theta}}\mathcal{L}_{SDS}(\phi, C_{\boldsymbol{\theta}}) = \mathbb{E}_{t,\epsilon}\left[w(t)(\epsilon_{\phi}(x_t; x_{lr}, t) - \epsilon)\frac{\partial x}{\partial \boldsymbol{\theta}}\right], \tag{2}$$

where $\phi$ is the pretrained image super-resolution diffusion model, $x_t$ is $x_0$ add noise $\epsilon$ at different diffusion timestep $t$, and $w(t)$ is a weight function of different noise levels.

Furthermore, to maintain the structural consistency and to prevent color shifts occasionally caused by diffusion model [32], the sub-pixel constraint $\mathcal{L}_{MSE}$ is also taken into consideration as a regularizer. The rendered high-resolution image $x_0$ is downsampled to align with its corresponding low-resolution image $x_{lr}$, which is formulated as follows:

$$\mathcal{L}_{MSE} = ||Downsample(x_0) - x_{lr}||. \tag{3}$$

In conclusion, the high-resolution 3DGS is joint optimized by $\mathcal{L}_{MSE} + \lambda\mathcal{L}_{SDS}$.

### 3.2 Gaussian Densification with SDS Constraint

In 3DGS, synthesizing high-quality novel views depends significantly on the representation capacity of Gaussian primitives. In particular, achieving accurate high-resolution rendering requires denser Gaussian primitives [33]. In our study, denser Gaussian primitives are produced by optimizing the high-resolution 3DGS with SDS. However, we observe that the direct application of SDS introduces undesirable and redundant Gaussian primitives during the densification process. We hypothesize that this issue arises from the inherent randomness of generative priors, as random noise and diffusion timesteps are sampled in the SDS process.

Referring to the training strategy of the diffusion model, SDS aims to optimize rendered high-resolution images to closely match the high-resolution distribution conditioned on its low-resolution counterparts by leveraging the data noising process. Nonetheless, during data noising, diffusion

4

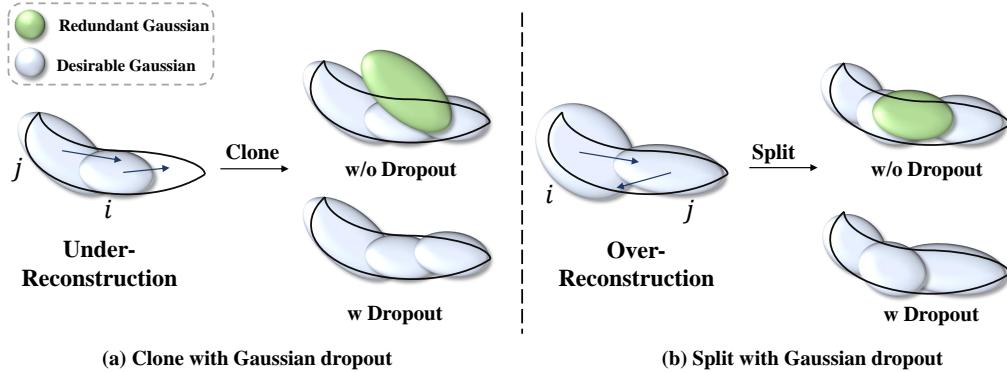**(a) Clone with Gaussian dropout**　　　　　　　　**(b) Split with Gaussian dropout**

Figure 3: Illustration of Gaussian Dropout during the densification process. When a small-scale object (depicted by the black outline) is insufficiently covered (under-reconstructed) or is represented by overly large splats (over-reconstructed), cloning or splitting is performed. In the top row (without dropout), a redundant Gaussian primitive (shown in green) is generated during densification. In the bottom row (with dropout), the redundant Gaussian primitive is randomly discarded.

timesteps are randomly sampled, resulting in varying gradient values with significant variance across different iterations. As described in Sec. 3.1, the Gaussian primitives with gradients exceeding a default threshold are transformed into two Gaussian primitives during densification. Consequently, the substantial variance in gradient values introduced by SDS leads to the generation of redundant Gaussian primitives.

Specifically, Fig. 2 (a) visualizes the gradient values across different diffusion timesteps $t$. As diffusion timesteps are randomly sampled in each training iteration, the significant variance of gradient values persists throughout the training process. Furthermore, we also visualize the variation of gradient values for a specific view during training under different constraints. The left figure in Fig. 2 (b) shows the gradient values under $\mathcal{L}_{SDS}$, which presents a large variance compared to the right figure under $\mathcal{L}_{MSE}$. Notably, the original 3DGS [2] employs $\mathcal{L}_{MSE}$ as the optimization constraint, which exhibits small variance across iterations and is well-suited to the default threshold strategy. In contrast, in our study, the high-variance gradient values brought by SDS, when subjected to the default threshold, can lead to the generation of redundant Gaussian primitives.

### 3.3 Stochastic Disturbance Reduction

To mitigate the aforementioned problem, we propose two techniques to reduce stochastic disturbances introduced by SDS: shrinking the range of diffusion timestep with an annealing strategy, and randomly discarding redundant Gaussian primitives during densification.

**Diffusion Timestep Annealing.** As a class of score-based generative models [34, 35, 36, 37], diffusion models involve a data noising and denoising process according to a predefined schedule over a fixed number of timesteps. Analogous to the training strategy of DDPM [34], the vanilla SDS randomly samples diffusion timestep $t$ from a uniform distribution (i.e., $t \sim \mathcal{U}(1, T)$) throughout the 3D model optimization. As described in Sec. 3.2, random sampling of diffusion timestep $t$ in SDS leads to redundant Gaussian primitives during the densification process. Therefore, we revise the timestep sampling range in SDS with an annealing strategy to reduce stochastic disturbances.

Particularly, the vanilla timestep sampling strategy of SDS involves sampling $t$ from a fixed range $[1, T]$ during each data noising step. In our approach, we refine it by employing an annealing strategy to progressively shrink the lower bound of the diffusion timestep sampling range. Specifically, for the current iteration $i$, the sampling range is adjusted to $[LB(i), T]$, where the lower bound $LB(i)$ is calculated as follows:

$$LB(i) = T - \frac{i}{N}. \tag{4}$$

5

In this equation, $T$ represents the upper bound, and $N$ denotes the annealing interval. Consequently, the diffusion timestep $t$ is sampled from the interval $[LB(i), T]$, i.e., $t \sim \mathcal{U}(LB(i), T)$, during each data noising step in the SDS process.

**Gaussian Dropout.**    In addition to reducing stochastic disturbances by shrinking the diffusion timestep sampling range, we directly discard undesirable and redundant Gaussian primitives using Gaussian Dropout. Specifically, as depicted in Fig. 3 (a), when considering the cloning of Gaussian primitive $i$ to fill the empty area (referred to as the "under-reconstruction" region), the nearby Gaussian primitive $j$, which should not be cloned, may exhibit large gradients due to disturbances from the SDS loss. This can lead to the generation of redundant Gaussian primitives. Therefore, to mitigate this issue, we employ Gaussian Dropout to discard the cloning or splitting of certain Gaussian primitives.

In detail, given a set of Gaussian primitives $G = \{g_0, g_1, ..., g_n\} \in \mathbb{R}^n$, where all gradients exceed the default threshold $\tau_{pos}$, we first generate a mask $M$ randomly with a certain probability $p$ and then use the mask to select a subset $G^{'} = \{g_0, g_2, ..., g_{n-2}, g_n\} \in \mathbb{R}^k (k < n)$ of $G$. The Gaussian primitives in subset $G^{'}$ will be split or cloned during densification, while the other Gaussian primitives will be dropped out and remain unchanged. Thus, the denser set $\hat{G} \in \mathbb{R}^{n+k}$ after densification can be formulated as:

$$\hat{G} = \mathcal{D}(G^{'}) + (G - G^{'}), \ \text{where } G^{'} = G \cdot M(p), M(p) = \begin{cases} 0 & rand(G) < p \\ 1 & else \end{cases}, \quad (5)$$

where $\mathcal{D}$ means the densification step.

## 4    Experiments

In this section, we present a comprehensive set of qualitative and quantitative evaluations aimed to verify the effectiveness of our proposed GaussianSR. Additionally, we conduct ablation studies to systematically evaluate the impact and effectiveness of each individual component.

### 4.1    Datasets and Metrics

**Blender Dataset [1].**    Blender Dataset is a Realistic Synthetic $360°$ Dataset that contains 8 detailed synthetic objects with a resolution of $800 \times 800$. We follow the same training and testing data split strategy as the original 3DGS [2]. For each scene, 100 images are used for training and 200 images are used for testing. The input resolution is set to $200 \times 200$, and we super-resolve this low-resolution 3DGS by a factor of 4. The downsampling method used is the same as the one provided in the official 3DGS code.

**Mip-NeRF 360 Dataset [38].**    Mip-NeRF 360 consists of 9 real-world scenes with 5 outdoors and 4 indoors. Each of them is composed of a complex central object or area with a detailed background. Following the previous setup, we use $7/8$ of the images for training and take the remaining $1/8$ for testing in each scene. We downsample the training views by a factor 4 as low-resolution inputs to $\times 4$ HRNVS task.

**Deep Blending Dataset [39].**    Deep Blending is a real-world dataset. Following 3DGS [2], we select two scenes of Deep Blending to evaluate our method. We use 1/8 of all views for testing and the rest for training. We downsample the training views by a factor 4 as low-resolution inputs to $\times 4$ HRNVS task.

**Metrics.**    The quality of view synthesis is assessed relative to the ground truth from the same pose, employing four metrics: Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity Index Measure (SSIM) [40], LPIPS (VGG) [41] and Frames Per Second (FPS).

### 4.2    Implementation Details

We implement our method based on the open-source 3DGS code. Training consists of 30k iterations for indoor scenes and 10k iterations for other scenes. As for the off-the-shelf 2D super-resolution diffusion model, we opt for StableSR [42] as our backbone. For the annealing interval $N$ in Eq. 4,

Table 1: Quantitative comparison for HRNVS ($\times 4$) with previous works on Blender, Mip-NeRF 360, and Deep Blending Dataset.

| Method | Blender | | | | Mip-NeRF 360 | | | | Deep Blending | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | FPS↑ | PSNR↑ | SSIM↑ | LPIPS↓ | FPS↑ | PSNR↑ | SSIM↑ | LPIPS↓ | FPS↑ |
| 3DGS [2] | 21.78 | 0.868 | 0.104 | **192** | 20.28 | 0.581 | 0.420 | **33** | 26.64 | 0.854 | 0.312 | **60** |
| StableSR [42] | 23.57 | 0.854 | 0.207 | <1 | 21.83 | 0.467 | 0.383 | <1 | 23.93 | 0.708 | 0.325 | <1 |
| Bicubic | 27.23 | 0.911 | 0.115 | 93 | 25.14 | 0.618 | 0.406 | 27 | 28.01 | 0.864 | 0.330 | 40 |
| NeRF-SR [3] | 27.81 | 0.920 | 0.097 | <1 | – | – | – | – | – | – | – | – |
| Ours | **28.37** | **0.924** | **0.087** | 192 | **25.60** | **0.663** | **0.368** | 33 | **28.28** | **0.873** | **0.307** | 60 |



Ground Truth      StableSR      3DGS      Bicubic      NeRF-SR      Ours
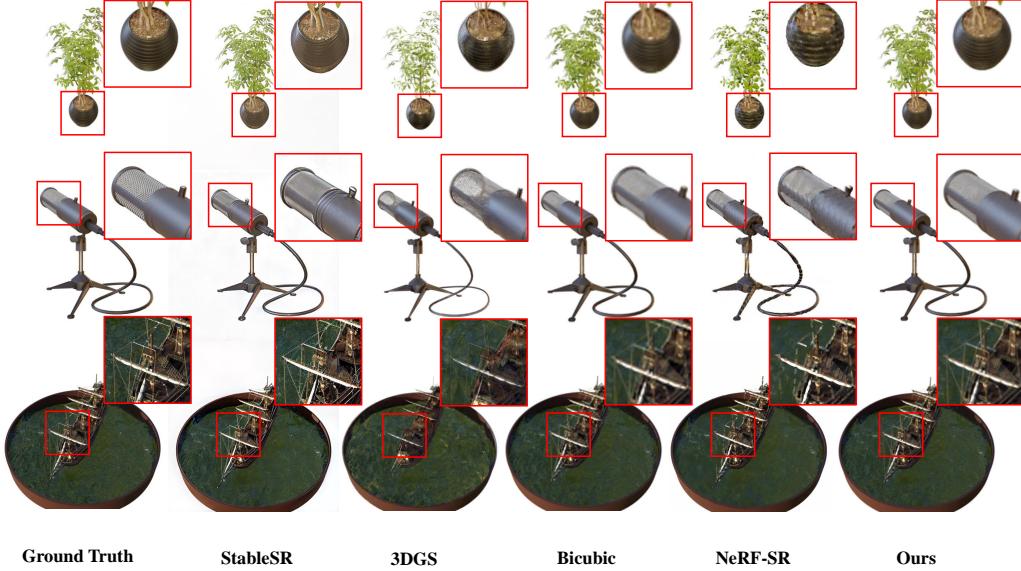
Figure 4: Qualitative comparison of the HRNVS ($\times 4$) on Blender dataset. Our method shows clearer details than 3DGS [2], Bicubic, NeRF-SR [3] and StableSR [42].

we shrink the sampling range of diffusion timestep every 100 iterations. The dropout probability $p$ of 0.7 is set during the Gaussian Dropout process. Additionally, bilinear interpolation is employed to downsample the rendered high-resolution images for $\mathcal{L}_{MSE}$ in Eq. 3, and $\lambda$ is set to be 0.001 during training. We perform experiments using a NVIDIA A100 GPU. To save space, please refer to our supplementary materials for more details.

### 4.3 Quantitative and Qualitative Comparisons

To demonstrate the effectiveness of our method, we compare it against several prior approaches, including vanilla 3DGS [2], bicubic interpolation, NeRF-SR [3] and StableSR [2]. For vanilla 3DGS baseline, we train 3DGS [2] using low-resolution input views and then render it at high resolution. Bicubic interpolation is applied to low-resolution renderings from the baseline 3DGS, providing a standard upsampling method. Regarding NeRF-SR [3], we directly run the source code to obtain qualitative and quantitative results. However, due to training instabilities encountered with Mip-NeRF 360 [38] and Deep Blending Dataset [39], we reproduce the results of NeRF-SR only on the Blender dataset [1]. For StableSR [42], we super-resolve each low-resolution view rendered from the baseline 3DGS using StableSR. Notably, since the 2D diffusion model we adopted (*i.e.*, StableSR [42]) is primarily trained on $4\times$ super-resolution data, we also primarily validate our GaussianSR on $\times 4$ HRNVS.

**Quantitative Evaluation.** Tab. 1 presents quantitative comparison results for $\times 4$ HRNVS tasks on the Blender dataset [1], the Mip-NeRF 360 dataset [38], and the Deep Blending dataset [39]. Our proposed GaussianSR outperforms previous state-of-the-art methods significantly in terms of PSNR, SSIM, and LPIPS metrics, while also requiring less rendering time. This indicates that GaussianSR excels in synthesizing high-resolution views with both superior quality and efficiency. Furthermore,
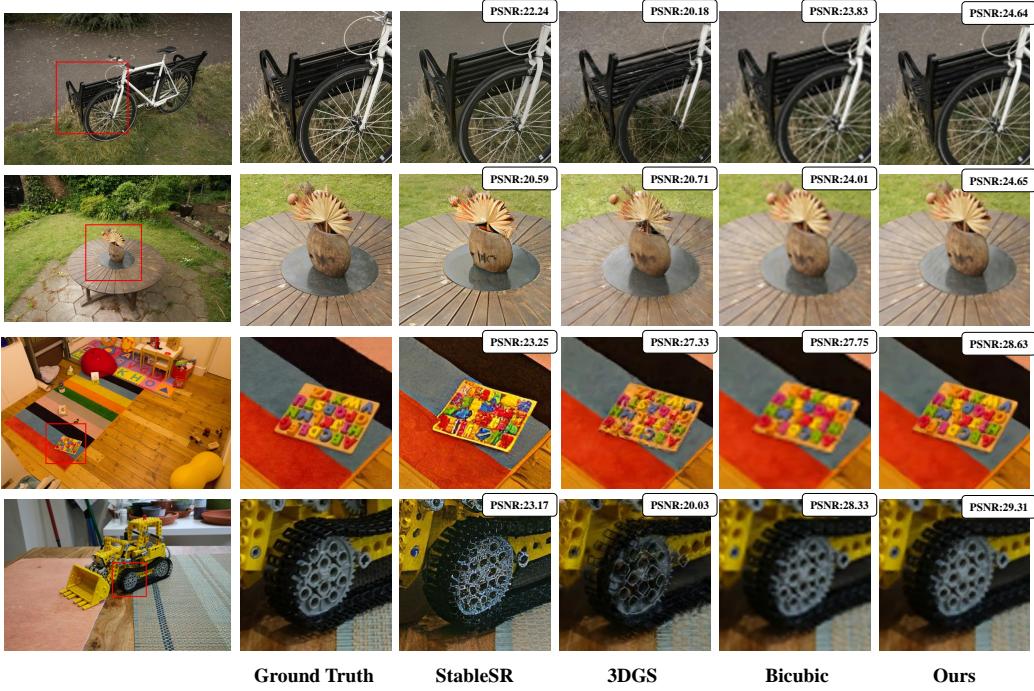
Figure 5: Qualitative comparison of our method with vanilla 3DGS, bicubic interpolation, and StableSR on Mip-NeRF 360 and Deep Blending Dataset for the HRNVS (×4). The results are the zoom-in version of the red box region and the PNSR value for the current view is presented in the top right corner. Our method presents higher quality and clearer details than others.

Table 2: Ablation studies on Mip-NeRF 360 and Deep Blending dataset for HRNVS (×4).

| Method | Mip-NeRF 360 | | | Deep Blending | | |
|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| Baseline (3DGS) | 20.28 | 0.581 | 0.420 | 26.64 | 0.854 | 0.312 |
| + $\mathcal{L}_{MSE}$ | 24.95 | 0.633 | 0.371 | 27.85 | 0.867 | 0.316 |
| + $\mathcal{L}_{SDS}$ (w dropout) | 25.36 | 0.631 | 0.369 | 28.18 | 0.874 | 0.311 |
| + Diffusion timestep annealing | **25.60** | **0.663** | **0.368** | **28.28** | **0.873** | **0.307** |

our method demonstrates the capability to generate detailed high-resolution novel views solely from low-resolution inputs, across synthetic as well as real-world datasets.

**Qualitative Evaluation.** Fig. 4 presents the qualitative results on the Blender dataset [1], while Fig. 5 shows the qualitative results on the Mip-NeRF 360 [38] and Deep Blending dataset [39]. GaussianSR consistently exhibits high-quality visual results across various scenarios, encompassing indoor and outdoor scenes. In contrast, the baseline model, 3DGS [2], suffers from needle-like artifacts due to the out-distribution rendering, whereas bicubic interpolation yields blurring artifacts by directly interpolating low-resolution views. Results super-resolved by StableSR [42] directly appear coarse and suffer from color shifts. Across both synthesis and real-world datasets, our GaussianSR produces superior visual results characterized by clearer edges and sharper details compared to alternative methods.

## 4.4 Ablation Studies

In Tab. 2, we perform ablation experiments on the components proposed in GaussianSR. Initially, we train the baseline 3DGS model with low-resolution inputs and subsequently render the high-resolution views directly. The effectiveness of each proposed component in GaussianSR is evaluated by gradually incorporating them into the model. In the second row of Table 2, we optimize the high-resolution 3DGS solely using $\mathcal{L}_{MSE}$. Subsequently, in the third row, we incorporate $\mathcal{L}_{SDS}$ and

8

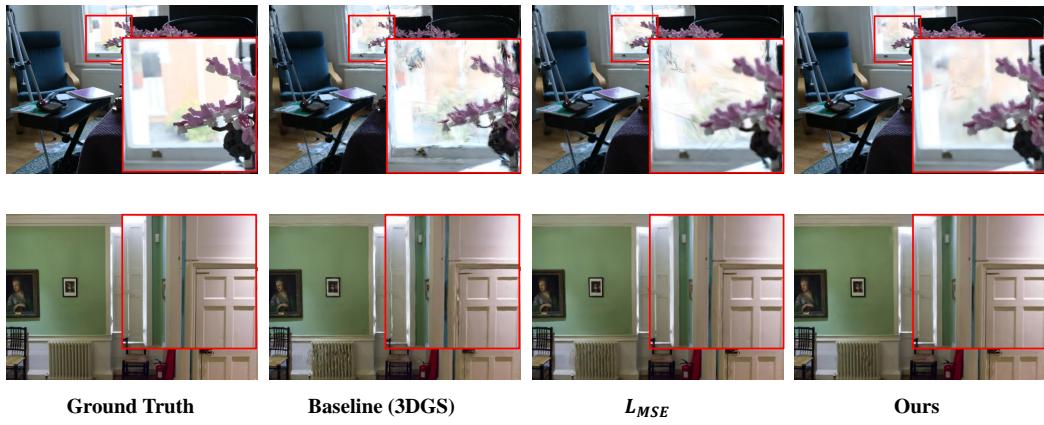| Ground Truth | Baseline (3DGS) | $L_{MSE}$ | Ours |

Figure 6: Qualitative evaluation for ablation studies. The third column means the results of high-resolution 3DGS that are optimized with $\mathcal{L}_{MSE}$ only. The last column presents the results of our full model. This demonstrates that our $\mathcal{L}_{SDS}$ with Gaussian Dropout and diffusion timestep annealing further yield clearer details.
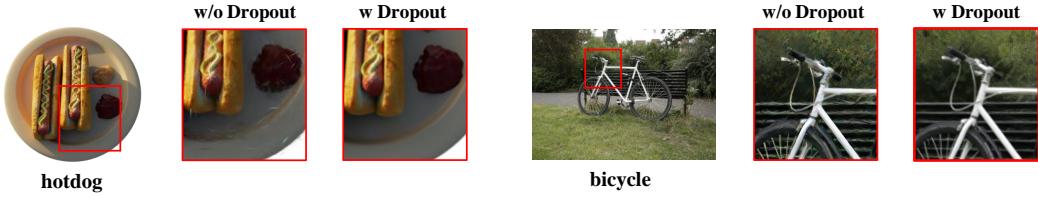


Figure 7: Ablation study on Gaussian Dropout.

Gaussian Dropout based on the second row. The effectiveness of the diffusion timestep annealing strategy is evaluated in the last row. Quantitative results unequivocally demonstrate that our proposed GaussianSR substantially enhances the quality of high-resolution novel views synthesized solely from low-resolution inputs. Additionally, we visualize the renderings of high-resolution novel views to assess the efficacy of our proposed components. As illustrated in Figure 6, our method effectively mitigates the presence of artifacts that may be present in 3DGS renderings. Furthermore, Figure 7 showcases results with and without Gaussian Dropout, revealing a notable reduction in redundant Gaussian primitives facilitated by Gaussian Dropout.

## 5 Conclusion

In this paper, we propose GaussianSR, an innovative method for synthesizing high-resolution novel views from low-resolution inputs. Our approach is grounded in 3D Gaussian Splatting (3DGS), which offers faster rendering speed. To address the challenge of limited high-resolution data, we employ Score Distillation Sampling (SDS) to distill generative priors of 2D super-resolution diffusion models. However, the direct application of SDS can lead to redundant Gaussian primitives due to the inherent randomness of generative priors. To mitigate this issue, we propose two straightforward yet effective techniques to reduce stochastic disturbance introduced by SDS. Experimental results demonstrate that GaussianSR excels in synthesizing higher-quality high-resolution novel views.

**Limitation and Future Works.** While our method shows promising results in high-resolution novel view synthesis (HRNVS), there remain limitations to be improved in future work. Our reliance on priors distilled from 2D super-resolution models constrains our performance to the capabilities of the specific 2D models employed. Future improvements could involve distilling priors from multiple 2D super-resolution models trained on diverse datasets, potentially enhancing performance and generalization.

# References

[1] 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. *Communications of the ACM*, 65(1):99–106, 2021.

[2] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), July 2023.

[3] Chen Wang, Xian Wu, Yuan-Chen Guo, Song-Hai Zhang, Yu-Wing Tai, and Shi-Min Hu. Nerf-sr: High quality neural radiance fields using supersampling. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 6445–6454, 2022.

[4] Yuqi Han, Tao Yu, Xiaohang Yu, Yuwang Wang, and Qionghai Dai. Super-nerf: View-consistent detail generation for nerf super-resolution. *arXiv preprint arXiv:2304.13518*, 2023.

[5] Xudong Huang, Wei Li, Jie Hu, Hanting Chen, and Yunhe Wang. Refsr-nerf: Towards high fidelity and super resolution view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8244–8253, 2023.

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

[7] Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Xinggang Wang. 4d gaussian splatting for real-time dynamic scene rendering. *arXiv preprint arXiv:2310.08528*, 2023.

[8] Zeyu Yang, Hongye Yang, Zijie Pan, Xiatian Zhu, and Li Zhang. Real-time photorealistic dynamic scene representation and rendering with 4d gaussian splatting. *arXiv preprint arXiv:2310.10642*, 2023.

[9] Zhicheng Lu, Xiang Guo, Le Hui, Tianrui Chen, Min Yang, Xiao Tang, Feng Zhu, and Yuchao Dai. 3d geometry-aware deformable gaussian splatting for dynamic view synthesis. *arXiv preprint arXiv:2404.06270*, 2024.

[10] Yi-Hua Huang, Yang-Tian Sun, Ziyi Yang, Xiaoyang Lyu, Yan-Pei Cao, and Xiaojuan Qi. Sc-gs: Sparse-controlled gaussian splatting for editable dynamic scenes. *arXiv preprint arXiv:2312.14937*, 2023.

[11] Yang Fu, Sifei Liu, Amey Kulkarni, Jan Kautz, Alexei A Efros, and Xiaolong Wang. Colmap-free 3d gaussian splatting. *arXiv preprint arXiv:2312.07504*, 2023.

[12] Zhiwen Fan, Wenyan Cong, Kairun Wen, Kevin Wang, Jian Zhang, Xinghao Ding, Danfei Xu, Boris Ivanovic, Marco Pavone, Georgios Pavlakos, et al. Instantsplat: Unbounded sparse-view pose-free gaussian splatting in 40 seconds. *arXiv preprint arXiv:2403.20309*, 2024.

[13] Zehao Zhu, Zhiwen Fan, Yifan Jiang, and Zhangyang Wang. Fsgs: Real-time few-shot view synthesis using gaussian splatting. *arXiv preprint arXiv:2312.00451*, 2023.

[14] Jaeyoung Chung, Jeongtaek Oh, and Kyoung Mu Lee. Depth-regularized optimization for 3d gaussian splatting in few-shot images. *arXiv preprint arXiv:2311.13398*, 2023.

[15] Jiahe Li, Jiawei Zhang, Xiao Bai, Jin Zheng, Xin Ning, Jun Zhou, and Lin Gu. Dngaussian: Optimizing sparse-view 3d gaussian radiance fields with global-local depth normalization. *arXiv preprint arXiv:2403.06912*, 2024.

[16] Xiaoyu Zhou, Zhiwei Lin, Xiaojun Shan, Yongtao Wang, Deqing Sun, and Ming-Hsuan Yang. Drivinggaussian: Composite gaussian splatting for surrounding dynamic autonomous driving scenes. *arXiv preprint arXiv:2312.07920*, 2023.

[17] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023.

[18] Yixun Liang, Xin Yang, Jiantao Lin, Haodong Li, Xiaogang Xu, and Yingcong Chen. Lucid-dreamer: Towards high-fidelity text-to-3d generation via interval score matching. *arXiv preprint arXiv:2311.11284*, 2023.

[19] Jaeyoung Chung, Suyoung Lee, Hyeongjin Nam, Jaerin Lee, and Kyoung Mu Lee. Lu-ciddreamer: Domain-free generation of 3d gaussian splatting scenes. *arXiv preprint arXiv:2311.13384*, 2023.

[20] Junshu Tang, Tengfei Wang, Bo Zhang, Ting Zhang, Ran Yi, Lizhuang Ma, and Dong Chen. Make-it-3d: High-fidelity 3d creation from a single image with diffusion prior. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22819–22829, 2023.

[21] Letian Huang, Jiayang Bai, Jie Guo, and Yanwen Guo. Gs++: Error analyzing and optimal gaussian splatting. *arXiv preprint arXiv:2402.00752*, 2024.

[22] Samuel Rota Bulò, Lorenzo Porzi, and Peter Kontschieder. Revising densification in gaussian splatting. *arXiv preprint arXiv:2404.06109*, 2024.

[23] Ziyi Yang, Xinyu Gao, Yangtian Sun, Yihua Huang, Xiaoyang Lyu, Wen Zhou, Shaohui Jiao, Xiaojuan Qi, and Xiaogang Jin. Spec-gaussian: Anisotropic view-dependent appearance for 3d gaussian splatting. *arXiv preprint arXiv:2402.15870*, 2024.

[24] Kai Cheng, Xiaoxiao Long, Kaizhi Yang, Yao Yao, Wei Yin, Yuexin Ma, Wenping Wang, and Xuejin Chen. Gaussianpro: 3d gaussian splatting with progressive propagation. *arXiv preprint arXiv:2402.14650*, 2024.

[25] Xiang Feng, Yongbo He, Yubo Wang, Yan Yang, Zhenzhong Kuang, Yu Jun, Jianping Fan, et al. Srgs: Super-resolution 3d gaussian splatting. *arXiv preprint arXiv:2404.10318*, 2024.

[26] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016.

[27] Tom McReynolds and David Blythe. *Advanced graphics programming using OpenGL*. Elsevier, 2005.

[28] Aditya Sanghi, Hang Chu, Joseph G Lambourne, Ye Wang, Chin-Yi Cheng, Marco Fumero, and Kamal Rahimi Malekshan. Clip-forge: Towards zero-shot text-to-shape generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18603–18613, 2022.

[29] Aditya Sanghi, Rao Fu, Vivian Liu, Karl DD Willis, Hooman Shayani, Amir H Khasahmadi, Srinath Sridhar, and Daniel Ritchie. Clip-sculptor: Zero-shot generation of high-fidelity and diverse shapes from natural language. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18339–18348, 2023.

[30] Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge Belongie, and Bharath Hariharan. Pointflow: 3d point cloud generation with continuous normalizing flows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4541–4550, 2019.

[31] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015.

[32] Jooyoung Choi, Jungbeom Lee, Chaehun Shin, Sungwon Kim, Hyunwoo Kim, and Sungroh Yoon. Perception prioritized training of diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11472–11481, 2022.

[33] Zhiwen Yan, Weng Fei Low, Yu Chen, and Gim Hee Lee. Multi-scale 3d gaussian splatting for anti-aliased rendering. *arXiv preprint arXiv:2311.17089*, 2023.

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

[35] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.

[36] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.

[37] 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.

[38] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022.

[39] Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel Brostow. Deep blending for free-viewpoint image-based rendering. *ACM Transactions on Graphics (ToG)*, 37(6):1–15, 2018.

[40] Zhou Wang, Eero P Simoncelli, and Alan C Bovik. Multiscale structural similarity for image quality assessment. In *The Thrity-Seventh Asilomar Conference on Signals, Systems & Computers, 2003*, volume 2, pages 1398–1402. Ieee, 2003.

[41] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018.

[42] Jianyi Wang, Zongsheng Yue, Shangchen Zhou, Kelvin CK Chan, and Chen Change Loy. Exploiting diffusion prior for real-world image super-resolution. In *arXiv preprint arXiv:2305.07015*, 2023.

# A Discussion of Hyperparameters

In this section, we discuss the hyperparameters selected in our method, including the weight-balancing parameter $\lambda$ of $\mathcal{L}_{MSE}$ and $\mathcal{L}_{SDS}$ in Sec. A.1, the annealing interval $N$ in Sec. A.2 and Gaussian Dropout probability $p$ in Sec. A.3.

## A.1 Weight-Balancing Parameter $\lambda$

To alleviate the shortage of data, we propose to leverage off-the-shelf 2D diffusion priors distilled by $\mathcal{L}_{SDS}$. Meanwhile, to maintain the consistency of low-resolution views and to prevent color shifts occasionally caused by diffusion model, we take $\mathcal{L}_{MSE}$ into consideration as a regularizer. Then, the high-resolution 3DGS is optimized by $\mathcal{L}_{MSE} + \lambda\mathcal{L}_{SDS}$, achieving high-resolution novel view synthesis. In this section, we make an analysis for the wight-balancing parameter $\lambda$. We randomly select three scenes form Blender dataset, Mip-NeRF 360 dataset and Deep Blending dataset to evaluate the performance under different $\lambda$. Tab. 3 presents that the best performance across three views in terms of PSNR and SSIM can be attained when $\lambda = 0.001$, whereas GaussianSR performs best in LPIPS when $\lambda$ is set to other values in "stump" and "playroom". After evaluating all aspects, we chose $\lambda = 0.001$ for training.

## A.2 Annealing Interval $N$

In order to reduce the randomness brought by generative priors, we shrink the diffusion timestep sampling range by an annealing strategy. In this section, we conduct an ablation of the annealing interval $N$ in Eq. 4. We evaluate the qualitative results under different $N$ on three scenes randomly selected from Mip-NeRF 360 [38] and Deep Blending [39] dataset. As shown in Tab. 4, GaussianSR achieves higher PSNR in the three scenes when $N$ is set to 100. Therefore, we shrink the diffusion timestep range every 100 iterations during training.

## A.3 Gaussian Dropout Probability $p$

As described in Eq. 5, we utilize the certain probability $p$ to generate a mask for suppress the cloning and splitting of some Gaussian primitives. Therefore, the performance is heavily related to the probability $p$. To chosen the $p$ with higher performance, we conduct an ablation under different dropout probabilities $p$ on the Blender dataset [1]. Referring to Fig. 8, GaussianSR demonstrates the best performance in terms of PSNR when $p = 0.7$, whereas it performs best in LPIPS when $p = 0.9$. Taking all aspects into consideration, $p = 0.7$ is chosen in the training process.

# B Additional Results

In this section, we provide more qualitative and quantitative results on Blender dataset [1], Mip-NeRF 360 dataset [38], and Deep Blending dataset [39]. For Blender dataset, Tab. 5 presents per-scene metrics for $\times4$ HRNVS. For each scene, we calculate the arithmetic mean of each metric averaged over all test views. More qualitative comparison on Blender dataset against leading methods is shown in Fig. 9. And per-scene metrics for $\times4$ HRNVS on Mip-NeRF 360 dataset are shown in Tab. 6, which demonstrates that our GaussianSR has the ability to synthesize higher-quality high-resolution novel views in most scenes. Following 3DGS [2], we select a subset of Deep Blending dataset to evaluate our method, where "drjohnson" and "playroom" are chosen. And the per-scene metrics of "drjohnson" and "playroom" are compiled in Tab. 7. Furthermore, more qualitative evaluation on Mip-NeRF 360 and Deep Blending dataset are presented in Fig. 10, Fig. 11 and Fig. 12. We also provide the video of our results in the supplementary materials which can entirely show the strength of our method.

Table 3: Ablation studies for weight-balancing parameter $\lambda$ on Blender dataset, Mip-NeRF 360 dataset and Deep Blending dataset.

| $\lambda$ | *ficus* | | | *stump* | | | *playroom* | | |
|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| 1 | 24.47 | 0.927 | 0.071 | 24.31 | 0.571 | 0.402 | 25.10 | 0.834 | 0.362 |
| 0.1 | 24.37 | 0.927 | 0.071 | 19.70 | 0.372 | 0.621 | 28.58 | 0.878 | 0.308 |
| 0.01 | 24.35 | 0.927 | 0.071 | 24.33 | 0.571 | 0.402 | 26.83 | 0.853 | 0.353 |
| 0.001 | 29.19 | 0.952 | 0.052 | 24.38 | 0.574 | 0.408 | 28.76 | 0.881 | 0.309 |
| 0.0001 | 28.73 | 0.948 | 0.056 | 24.28 | 0.571 | 0.401 | 25.20 | 0.835 | 0.362 |

Table 4: Ablation studies for annealing interval $N$ on Mip-NeRF 360 dataset and Deep Blending dataset .

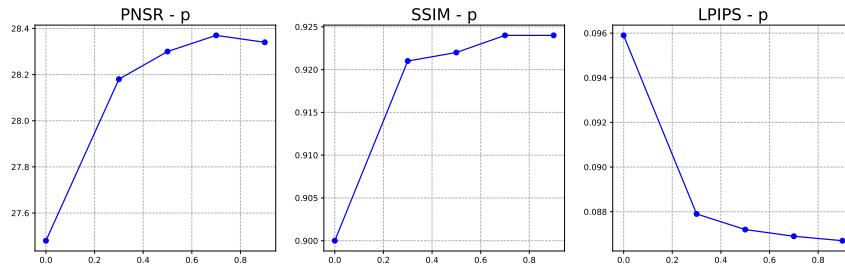| $N$ | *kitchen* | | | *treehill* | | | *drjohnson* | | |
|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| 500 | 28.19 | 0.790 | 0.282 | 21.80 | 0.482 | 0.473 | 27.79 | 0.864 | 0.307 |
| 300 | 28.02 | 0.791 | 0.278 | 21.86 | 0.485 | 0.473 | 27.80 | 0.866 | 0.305 |
| 100 | 28.27 | 0.794 | 0.276 | 21.98 | 0.490 | 0.478 | 27.81 | 0.865 | 0.307 |



Figure 8: Ablation studies of Gaussian Dropout probability $p$ on Blender dataset.

Table 5: Quantitative evaluation for HRNVS (×4) on the Blender dataset. For each scene, we report the arithmetic mean of each metric averaged over all test views.

**PSNR**

| Method | chair | drums | ficus | hotdog | lego | materials | mic | ship | Average |
|---|---|---|---|---|---|---|---|---|---|
| 3DGS [2] | 24.14 | 19.33 | 21.72 | 25.97 | 20.46 | 20.27 | 21.44 | 20.90 | 21.78 |
| StableSR [42] | 24.69 | 21.12 | 23.45 | 26.68 | 22.64 | 22.98 | 24.88 | 22.17 | 23.58 |
| Bicubic | 28.81 | 23.35 | 27.32 | 31.67 | 27.19 | 26.07 | 27.90 | 25.51 | 27.23 |
| NeRF-SR [3] | 30.18 | 23.50 | 22.72 | 34.38 | 29.21 | 28.08 | 27.25 | 26.59 | 27.81 |
| Ours | 29.81 | 24.05 | 29.19 | 32.80 | 28.66 | 27.02 | 29.15 | 26.28 | 28.37 |

**SSIM**

| Method | chair | drums | ficus | hotdog | lego | materials | mic | ship | Average |
|---|---|---|---|---|---|---|---|---|---|
| 3DGS [2] | 0.886 | 0.852 | 0.916 | 0.920 | 0.821 | 0.859 | 0.914 | 0.778 | 0.868 |
| StableSR [42] | 0.870 | 0.859 | 0.893 | 0.890 | 0.808 | 0.856 | 0.916 | 0.740 | 0.854 |
| Bicubic | 0.918 | 0.892 | 0.934 | 0.951 | 0.902 | 0.916 | 0.947 | 0.824 | 0.911 |
| NeRF-SR [3] | 0.937 | 0.903 | 0.904 | 0.964 | 0.931 | 0.932 | 0.943 | 0.836 | 0.919 |
| Ours | 0.931 | 0.909 | 0.952 | 0.959 | 0.924 | 0.925 | 0.958 | 0.837 | 0.924 |

**LPIPS**

| Method | chair | drums | ficus | hotdog | lego | materials | mic | ship | Average |
|---|---|---|---|---|---|---|---|---|---|
| 3DGS [2] | 0.083 | 0.108 | 0.062 | 0.077 | 0.142 | 0.107 | 0.061 | 0.192 | 0.104 |
| StableSR [42] | 0.178 | 0.178 | 0.201 | 0.200 | 0.209 | 0.227 | 0.197 | 0.265 | 0.207 |
| Bicubic | 0.091 | 0.118 | 0.073 | 0.078 | 0.132 | 0.096 | 0.063 | 0.219 | 0.115 |
| NeRF-SR [3] | 0.068 | 0.108 | 0.100 | 0.053 | 0.090 | 0.076 | 0.078 | 0.198 | 0.097 |
| Ours | 0.074 | 0.095 | 0.052 | 0.061 | 0.106 | 0.083 | 0.050 | 0.175 | 0.087 |

Table 6: Quantitative evaluation for HRNVS (×4) on the Mip-NeRF 360 dataset. For each scene, we report the arithmetic mean of each metric averaged over all test views.

**PSNR**

| Method | bicycle | flowers | garden | stump | treehill | bonasi | counter | kitchen | room |
|---|---|---|---|---|---|---|---|---|---|
| 3DGS [2] | 18.48 | 15.80 | 19.73 | 18.65 | 15.91 | 23.18 | 23.52 | 20.66 | 26.58 |
| StableSR [42] | 20.61 | 18.28 | 20.77 | 21.83 | 20.23 | 23.32 | 23.33 | 22.43 | 24.61 |
| Bicubic | 22.30 | 20.43 | 23.27 | 24.20 | 22.05 | 29.43 | 27.38 | 27.40 | 29.78 |
| Ours | 22.63 | 20.61 | 23.72 | 24.38 | 21.98 | 30.23 | 28.02 | 28.27 | 30.58 |

**SSIM**

| Method | bicycle | flowers | garden | stump | treehill | bonasi | counter | kitchen | room |
|---|---|---|---|---|---|---|---|---|---|
| 3DGS [2] | 0.405 | 0.326 | 0.463 | 0.405 | 0.404 | 0.773 | 0.754 | 0.694 | 0.826 |
| StableSR [42] | 0.337 | 0.315 | 0.398 | 0.409 | 0.384 | 0.651 | 0.656 | 0.546 | 0.678 |
| Bicubic | 0.491 | 0.436 | 0.516 | 0.567 | 0.489 | 0.871 | 0.823 | 0.759 | 0.861 |
| Ours | 0.511 | 0.453 | 0.556 | 0.574 | 0.490 | 0.883 | 0.839 | 0.794 | 0.873 |

**LPIPS**

| Method | bicycle | flowers | garden | stump | treehill | bonasi | counter | kitchen | room |
|---|---|---|---|---|---|---|---|---|---|
| 3DGS [2] | 0.481 | 0.509 | 0.446 | 0.476 | 0.526 | 0.353 | 0.330 | 0.336 | 0.322 |
| StableSR [42] | 0.385 | 0.440 | 0.374 | 0.424 | 0.415 | 0.360 | 0.340 | 0.373 | 0.336 |
| Bicubic | 0.486 | 0.509 | 0.455 | 0.442 | 0.512 | 0.292 | 0.318 | 0.313 | 0.327 |
| Ours | 0.450 | 0.448 | 0.395 | 0.408 | 0.478 | 0.273 | 0.283 | 0.276 | 0.303 |

Table 7: Quantitative evaluation for HRNVS (×4) on the Deep Blending dataset. For each scene, we report the arithmetic mean of each metric averaged over all test views.

**PSNR**

| Method | drjohnson | playroom | Average |
|---|---|---|---|
| 3DGS [2] | 26.64 | 27.17 | 26.64 |
| StableSR [42] | 24.02 | 23.83 | 23.93 |
| Bicubic | 27.73 | 28.29 | 28.00 |
| Ours | 27.81 | 28.76 | 28.28 |

**SSIM**

| Method | drjohnson | playroom | Average |
|---|---|---|---|
| 3DGS [2] | 0.847 | 0.860 | 0.853 |
| StableSR [42] | 0.706 | 0.710 | 0.708 |
| Bicubic | 0.858 | 0.870 | 0.864 |
| Ours | 0.865 | 0.881 | 0.873 |

**LPIPS**

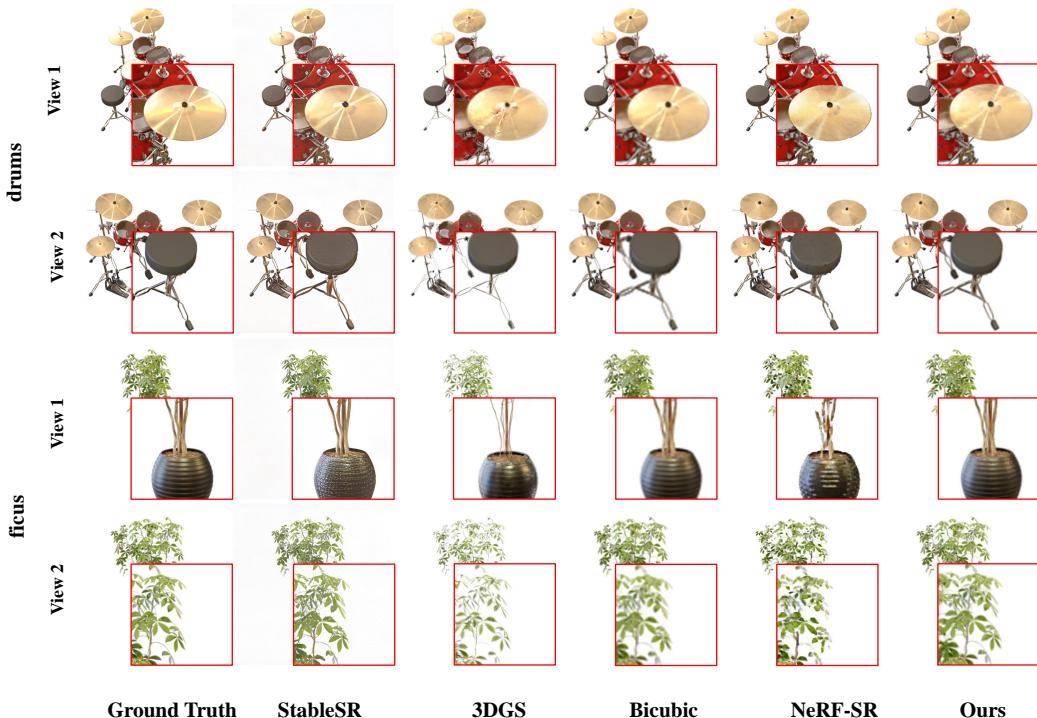| Method | drjohnson | playroom | Average |
|---|---|---|---|
| 3DGS [2] | 0.313 | 0.312 | 0.312 |
| StableSR [42] | 0.332 | 0.318 | 0.325 |
| Bicubic | 0.329 | 0.330 | 0.330 |
| Ours | 0.307 | 0.309 | 0.308 |



Figure 9: Qualitative comparison of our method with vanalia 3DGS, bicubic interpolation, NeRF-SR and StableSR on Blender dataset for the HRNVS (×4).
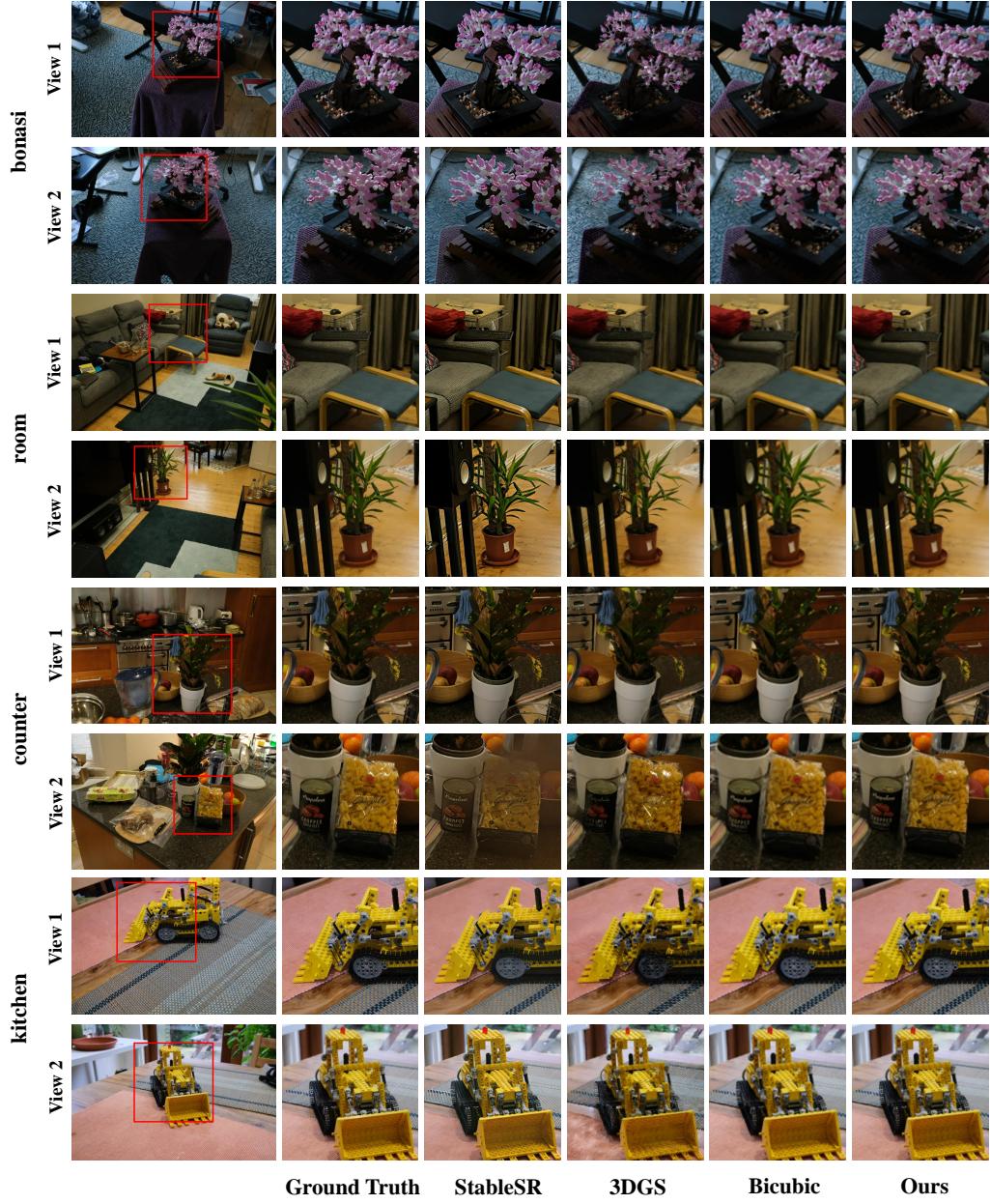
Figure 10: Qualitative comparison of our method with vanalia 3DGS, bicubic interpolation and StableSR in the indoor scenes of Mip-NeRF 360 dataset for the HRNVS (×4). The results are zoom-in version of the red box region.

Figure 11: Qualitative comparison of our method with vanalia 3DGS, bicubic interpolation and StableSR in the outdoor scenes of Mip-NeRF 360 dataset for the HRNVS (×4). The results are zoom-in version of the red box region.
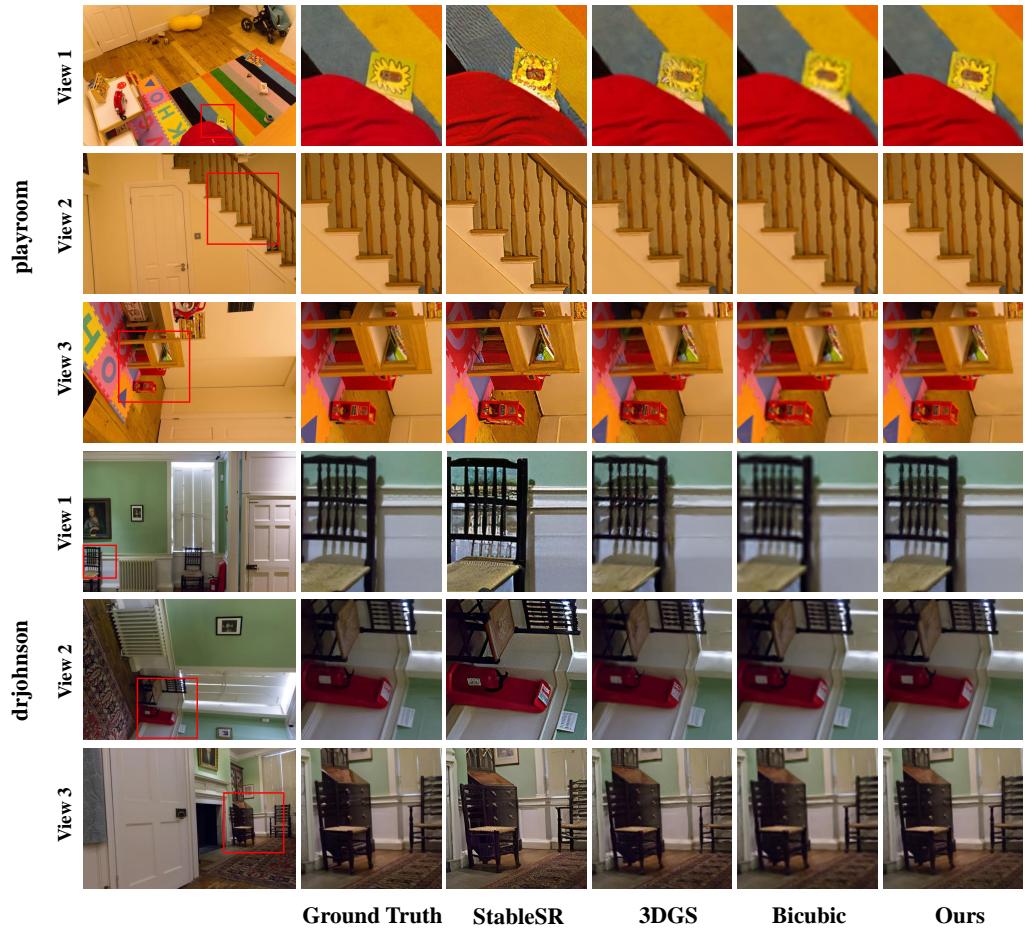
Figure 12: Qualitative comparison of our method with vanalia 3DGS, bicubic interpolation and StableSR on Deep Blending dataset for the HRNVS (×4). The results are zoom-in version of the red box region.