

A Lesson in Splats: Teacher-Guided Diffusion for 3D Gaussian Splats Generation with 2D Supervision

Chensheng Peng¹ Ido Sobol² Masayoshi Tomizuka¹ Kurt Keutzer¹ Chenfeng Xu¹ Or Litany^{2,3}

¹ UC Berkeley ² Technion ³ NVIDIA

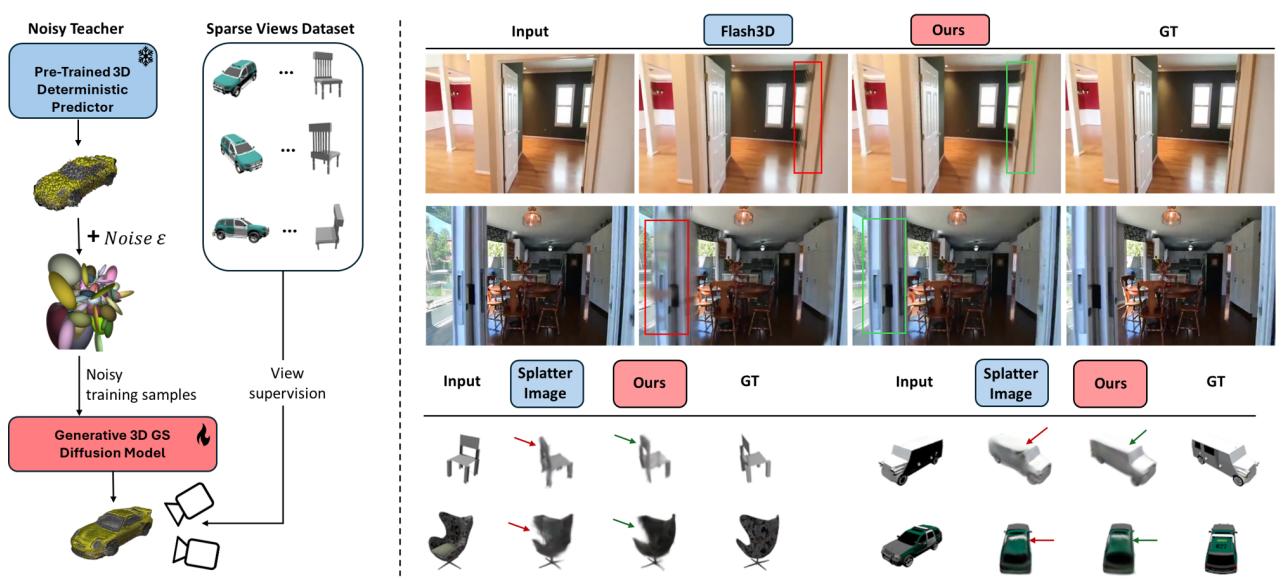


Figure 1. (Left) Standard diffusion training is typically constrained to same-modality supervision. We introduce a framework that breaks this barrier by decoupling the sources of noised samples and supervision. Leveraging imperfect predictions of a pretrained feedforward 3D reconstruction module, we train a denoiser that operates in 3D, while being supervised by sparse 2D views. (Right) When paired with two different noisy teachers, our diffusion model enhances reconstruction quality across both objects and scenes, particularly in regions where deterministic models struggle due to reconstruction ambiguity.

Abstract

We introduce a diffusion model for Gaussian Splats, Splat-Diffusion, to enable generation of three-dimensional structures from single images, addressing the ill-posed nature of lifting 2D inputs to 3D. Existing methods rely on deterministic, feed-forward predictions, which limit their ability to handle the inherent ambiguity of 3D inference from 2D data. Diffusion models have recently shown promise as powerful generative models for 3D data, including Gaussian splats; however, standard diffusion frameworks typically require the target signal and denoised signal to be in the same modality, which is challenging given the scarcity of 3D data. To overcome this, we propose a novel training strategy that decouples the denoised modality from the supervision modality. By using a deterministic model as a noisy teacher to create

the noised signal and transitioning from single-step to multi-step denoising supervised by an image rendering loss, our approach significantly enhances performance compared to the deterministic teacher. Additionally, our method is flexible, as it can learn from various 3D Gaussian Splat (3DGS) teachers with minimal adaptation; we demonstrate this by surpassing the performance of two different deterministic models as teachers, highlighting the potential generalizability of our framework. Our approach further incorporates a guidance mechanism to aggregate information from multiple views, enhancing reconstruction quality when more than one view is available. Experimental results on object-level and scene-level datasets demonstrate the effectiveness of our framework.

1. Introduction

3D reconstruction from single 2D images is essential for computer vision applications, such as augmented reality, robotics, and autonomous vehicles, which rely on inferring 3D structures from limited viewpoints. Lifting 2D images to 3D is an ill-posed problem, because different 3D shapes can produce identical 2D projections.

Current approaches for 3D reconstruction from single images can be categorized into two main types: deterministic predictions and generative models, each with distinct limitations that stem from the ill-posed nature of inferring 3D structures from 2D images.

A prevalent approach in 3D reconstruction is to use deterministic neural networks to map input images to 3D representations, such as Neural Radiance Fields (NeRF) [33], 3D Gaussian Splats (3DGS) [39], and triplanes [68]. By leveraging differentiable rendering techniques, these methods can train directly from 2D images, circumventing the need for large volumes of 3D data, which are relatively scarce. This is advantageous because annotated 3D data is often difficult or impractical to obtain, especially for real-world applications. However, despite ongoing performance improvements, deterministic models remain inherently limited by the ambiguity in the 2D-to-3D mapping. These models cannot fully capture the range of possible 3D structures that correspond to a single 2D image, leading to overly smooth or blurred outputs when supervised by appearance-based losses for a specific target. This limitation constrains their ability to generate diverse and accurate 3D reconstructions, particularly when multiple plausible 3D structures could correspond to the same 2D input.

In contrast, generative models, especially diffusion models [8, 15], have recently shown strong potential in representing 3D data. Diffusion models are trained to progressively denoise corrupted versions of 3D data to generate 3D outputs that are likely under the training set distribution, either by directly operating in the 3D space [1, 30, 34] or in a lower-dimensional latent space [39, 40, 51]. These models have proven powerful for generating realistic data across various modalities. However, diffusion models for 3D generation face a fundamental limitation due to their training process, in which the denoising model is trained on noisy samples using their clean counterparts as supervision. This requirement demands a substantial amount of 3D data, making these models difficult to scale to real-world applications where 3D data is limited. Consequently, while diffusion models are highly effective, their dependency on large 3D datasets restricts their applicability for 3D generation when starting from 2D image data.

Some attempts have been made to bypass these limitations by training 3D generative models using multi-view images [46]. These models aggregate information across multiple views, structuring predictions in 3D space. How-

ever, they still fall short in generation quality compared to deterministic methods and require multiple views for each target object, which is often impractical in real-world applications. Thus, although both deterministic and generative models have made strides in 3D reconstruction, the field lacks scalable, high-performance solutions that can infer 3D structures from single 2D images. Research into training 3D diffusion models from single-image inputs remains underexplored, highlighting an important gap that our work aims to address.

In this work, we propose a novel training strategy that fundamentally revises the principles of diffusion model training by decoupling the denoised modality from the supervision modality. To clarify the challenge, we note that although a denoiser can be trained in the 3D modality and supervised with 2D images, diffusion models require the noisy signal and the clean signal to remain in the same modality—here, in 3D. However, since our supervision comes from 2D images, the key question becomes: *how to create a noisy 3D signal suitable for denoising?*

Our solution leverages powerful deterministic 3D reconstruction methods as “noisy teachers”. While deterministic models can produce 3D representations that are imperfect, as manifested by blurry renderings, they nonetheless provide a useful starting point in 3D. This approach is advantageous because the goal in diffusion training is to apply noise to the signal. By introducing enough noise beyond a critical timestep t^* , the difference between the noisy 3D signal provided by the deterministic model and the true, clean 3D structure becomes minimal. This “sweet spot” in noise level, inspired by techniques like SDEdit [60], allows us to use the deterministic 3D output to create a suitable noisy input for diffusion training. At this stage, we can train a denoiser to predict the clean signal as is often done for guidance using Tweedie’s formula [10, 22, 38], and we can supervise this prediction using image-based rendering losses.

However, this alone is not sufficient because if the denoiser only learns from timesteps $t > t^*$, it is bound to producing blurry outputs. To overcome this, we introduce a second key innovation: a multi-step denoising strategy that replaces the traditional single-step denoising framework. Specifically, starting from a noise level $t > t^*$, our model performs iterative denoising, akin to its behavior during inference, progressively reducing noise over multiple steps until reaching the clean 3D structure. This multi-stage approach allows the model to generate sharp and detailed reconstructions, as each denoising step moves closer to the true 3D target. By supervising the final output with image-based rendering losses, we enable the model to learn to produce high-quality outputs while updates propagate across the entire denoising sequence. Consequently, our approach opens up the capability to directly train 3D diffusion models from 2D supervision.

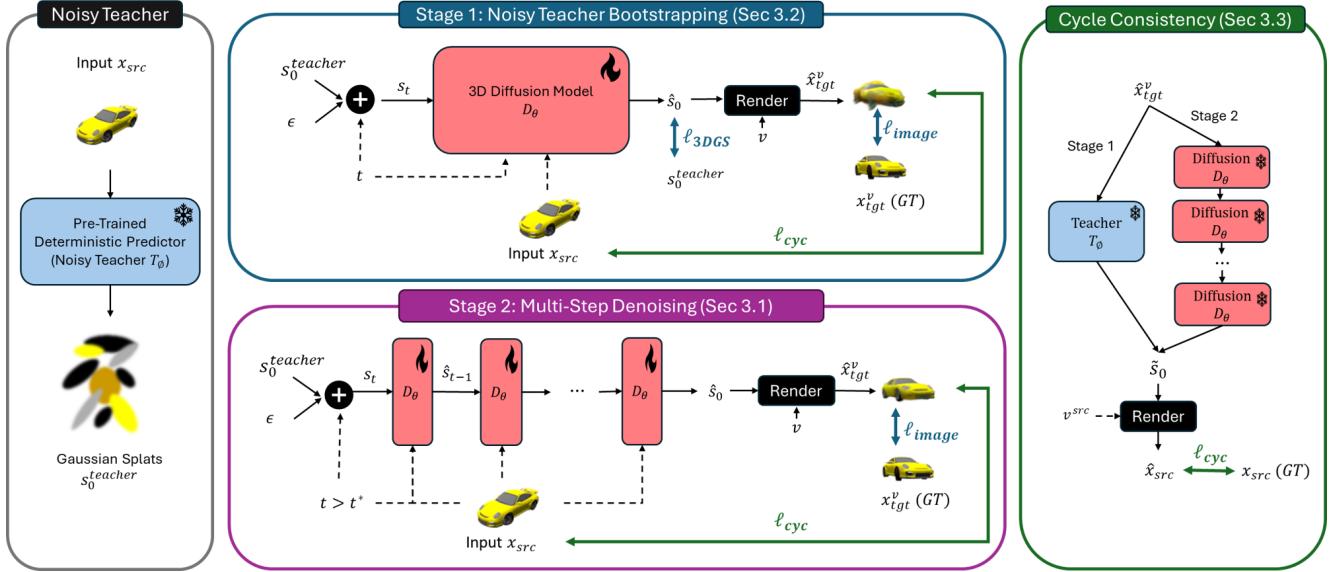


Figure 2. **Our proposed framework for noisy-teacher-guided training of a 3D Gaussian Splat (3DGS) diffusion model.** Using a pre-trained deterministic predictor network for 3DGS, which we refer to as the “noisy teacher” (left), in stage 1 (top) we lift sampled views to generate an imperfect 3DGS prediction, providing noisy samples and supervision for the diffusion denoiser in 3DGS with additional image supervision. In stage 2 (bottom), we decouple the noisy samples from supervision and instead use the noisy teacher to generate noisy samples at noise levels $t > t^*$, with a multi-step denoising strategy generating high-quality predictions to facilitate image-only supervision. Both stages incorporate cycle consistency regularization. See text for further details.

Notably, this strategy is extremely flexible and can utilize various 3D teacher models as long as the diffusion prediction can be supervised by them. In our experiments, we demonstrate this flexibility using two types of deterministic models: the popular Splatter Image [48] and the more advanced and recent Flash3D [47]. With these models, we train on single object and scene data, respectively. In both cases, our method significantly improves the performance of the base teacher model by $0.5 - 0.85$ PSNR. Additionally, our diffusion model facilitates the incorporation of additional views through guidance, further boosting performance compared to standard optimization.

2. Related Work

2.1. 3D Reconstruction from Sparse Views with Deterministic Models

Recent research has focused on generating 3D content from images using deterministic feed-forward models [17, 47, 48, 63]. Notably, these methods rely solely on posed 2D views for training, rather than requiring 3D data, making them scalable for in-the-wild training. While deterministic models are relatively simple to design and train, they struggle to capture the inherent variability of possible solutions in 3D reconstruction, often leading to blurry reconstructions in regions with large potential variability. In this work, we advocate for a generative 3D diffusion model to enable richer

and more complex representations. We use deterministic models [47, 48] as a starting point to generate noisy samples, which are then used to train our diffusion model.

2.2. 3D Generation with Diffusion Models

Diffusion models have shown impressive generative capabilities across various domains, leading to significant interest in applying them to 3D content generation.

Diffusion Models Trained Directly on 3D Data. One line of research focuses on designing diffusion models that directly operate in 3D space. These models have been developed for various 3D representations, including point clouds [31, 51, 65], meshes [1, 30], 3D Gaussian splats [34, 39], and neural fields [4, 7, 9, 35, 43]. While effective, these methods assume the availability of high-quality 3D datasets in the target representation, which are often scarce and lack the breadth of real-world diversity. This data scarcity limits the generalization and applicability of these models, particularly in in-the-wild scenarios.

Leveraging 2D Diffusion Models for 3D Content Creation. To address the scarcity of 3D data, recent works have explored leveraging 2D-trained diffusion models to create 3D content. A prominent technique in this line is Score Distillation Sampling (SDS), which “lifts” 2D score predictions to a shared 3D representation [14, 20, 25, 32, 36, 37, 53, 62]. However, a key challenge here is achieving view coherence, as 2D models only access the visible parts of an object, lead-

ing to potential issues such as the notorious Janus problem. To mitigate this, view-aware diffusion models, condition the generation of target views on one or more source views, incorporating relative camera transformations for enhanced coherence [5, 11, 16, 24, 28, 29, 42, 45, 54, 57–59].

3D Diffusion Models Supervised by 2D Images. Our work aligns with a relatively underexplored area focused on training diffusion models that operate in 3D space but are supervised only with 2D images. Traditionally, in diffusion models, the supervision signal is provided in the same modality as the noisy samples. Holodiffusion [19] introduced a method to train a 3D diffusion model for feature voxel grids using 2D supervision. To address the discrepancy between the noised samples and the noised target distribution, they apply an additional denoising pass, encouraging the model to learn both distributions simultaneously.

In contrast, our approach minimizes the distribution discrepancy between teacher-induced noised samples and (unavailable) target noise samples by focusing on large noise values and refining lower-noise predictions through a multi-step denoising process. Several approaches [2, 46, 49], denoise multi-view images using a denoiser *structured* to predict a 3D representation, which is then rendered into 2D views. However, these methods inherently rely on the bijectivity of multi-view and 3D representations, which only hold with a substantial number of images. Additionally, because the images are noised independently, they may not coherently represent the noisy 3D structure, potentially harming consistency. Our proposed method, in contrast, directly denoises within the 3D representation while using 2D views for supervision, addressing both data scarcity and view coherence by explicitly working in 3D space.

3. Method

Problem Formulation. We tackle the problem of training a 3D diffusion model, $D_\theta(s_t, t, x_{\text{src}})$, that maps N noisy 3D Gaussian Splats, $s_t \in \mathbb{R}^{N \times d}$ to their clean version s_0 . Each of the Gaussian Splats is of dimension d , representing properties such as center, covariance, opacity, and color. The model is conditioned on a single image x_{src} and uses $k \geq 1$ additional views of the same content for supervision, $\{x_{\text{tgt}}^v\}_{v=0}^{k-1}$, without access to 3D ground truth. We assume access to a pre-trained deterministic model $s_0^{\text{teacher}} = T_\phi(x)$, trained on the same sparse view data, that reconstructs 3D Gaussian Splats from a single image—or we can train such a model ourselves. Our method employs this trained model as a noisy teacher, generating noisy samples to train the diffusion model, which is supervised by the target image set $\{x_{\text{tgt}}^v\}_{v=0}^{k-1}$.

Overview. Our pipeline operates in two stages. First, we bootstrap the diffusion model by supervising it with the noisy teacher’s predictions (Section 3.2). We then proceed to fine-tune the diffusion model using multi-step denoising

and rendering losses (Section 3.1). Both stages are further equipped with a cycle consistency regularization described in Section 3.3. Although the bootstrapping stage precedes fine-tuning in the pipeline, we present it second in this manuscript to facilitate a smoother explanation of our core contributions. The model pipeline is depicted in Fig. 2.

3.1. Decoupling Noised Samples from Supervision with Multi-Step Denoising

Our approach to overcoming the aforementioned unimodality limitation of diffusion model training is to decouple the source for the noisy samples from the supervision. Specifically, in standard diffusion training, noise is added to the target ground truth sample, which is then fed to the denoiser for recovering the clean target. Here, we do not have access to true 3D target data; instead, we replace it with a 3D prediction from a pretrained deterministic model. As previously discussed this model is limited in its ability to generate the diverse plausible 3D structures often resulting in blurry and imprecise predictions, thus we consider it to be a “noisy teacher”. A key insight is that while the noisy teacher does not produce 3D Gaussian Splats (3DGS) that are sufficient as a standalone solution, they are useful as a starting point in our proposed framework. We further take inspiration from [60], which finds that with enough noise, the data distribution of two modalities can overlap. Based on this, we choose a timestep t^* such that for $t \geq t^*$, the noisy samples generated by the noisy teacher are likely to align with those that would have resulted from a forward noising process applied to the true, unknown ground truth 3DGS. Denoting these samples as

$$s_t = \sqrt{\alpha_t} s_0^{\text{teacher}} + \sqrt{1 - \alpha_t} \epsilon, \quad (1)$$

With the input image $x_{\text{src}} \sim p_{\text{data}}$ sampled from the image dataset and noise $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ —these notations are omitted for brevity throughout the manuscript. One might be tempted to train the denoiser using the standard training objective:

$$\mathbb{E}_{x_{\text{src}}, t > t^*, \epsilon} [\|s_0^{\text{teacher}} - D_\theta(s_t, t, x_{\text{src}})\|_2^2]. \quad (2)$$

However, a problem remains: the noise ϵ is the noise added to the noisy teacher, so predicting it would not help since it is not the noise from the unknown true target. Instead, we utilize the fact that the predicted 3DGS representation s can be differentiably rendered in arbitrary view directions v . Denoting this rendering operation as $\mathcal{R}(s, v)$, we can modify the training scheme to:

$$\mathbb{E}_{x_{\text{src}}, v \sim \mathcal{U}[k], t > t^*, \epsilon} [\|x_{\text{tgt}}^v - \mathcal{R}(D_\theta(s_t, t, x_{\text{src}}), v)\|_2^2]. \quad (3)$$

Yet, an issue still exists. By limiting our sample range of timesteps, we do not sample small noise levels, and as a result, the model cannot recover the fine details essential for successful reconstruction. Sampling smaller timesteps is not

ideal, as the model would then be trained on noisy samples from the incorrect distribution.

To address this, we revise the standard single-step denoising training and instead employ *multi-step denoising*, sequentially applying the model with the appropriate time-step conditioning until reaching the final clean 3D prediction, $\hat{s}_0 = D_\theta(\hat{s}_1, 1, x_{\text{src}}) \circ \dots \circ D_\theta(s_t, t, x_{\text{src}})$. Rendered towards a target view, the loss becomes:

$$\mathcal{L}_{\text{mlt-stp}} = \mathbb{E}_{x_{\text{src}}, v \sim \mathcal{U}[k], t > t^*, \epsilon} [\lambda_t \|x_{\text{tgt}}^v - \mathcal{R}(\hat{s}_0, v)\|_2^2], \quad (4)$$

where λ_t assigns a weight per denoising step. This multi-step denoising process mirrors the inference process but allows the network parameters to update. By **training** the model in this way, the 3D denoiser learns to handle 3D data directly, while still being supervised using widely available 2D datasets. Please refer to the implementation details 4.2 for a discussion regarding the computational efficiency of this unrolled optimization.

3.2. Noisy Teacher Bootstrapping

Training a 3D diffusion model directly using the multi-step denoising paradigm described above is computationally expensive. This is primarily due to the increased memory costs of maintaining gradients over multiple denoising steps in 3D space, which limits batch sizes and reduces efficiency. To address this, we propose avoiding this training approach from scratch by first bootstrapping our model using the noisy teacher.

Specifically, we generate noisy samples s_t from the noisy teacher, as shown in Equation 1, and supervise the generated 3DGS both directly in 3D:

$$\ell_{\text{3DGS}} = \|s_0^{\text{teacher}} - D_\theta(s_t, t, x_{\text{src}})\|^2, \quad (5)$$

and in 2D through the image rendered from the generated 3DGS:

$$\ell_{\text{image}} = \|x_{\text{tgt}}^v - \mathcal{R}(D_\theta(s_t, t, x_{\text{src}}), v)\|_2^2. \quad (6)$$

These losses are combined to form our overall bootstrapping objective:

$$\mathcal{L}_{\text{bootstrap}} = \mathbb{E}_{x_{\text{src}}, v \sim \mathcal{U}[k], t \sim \mathcal{U}[T], \epsilon} [\ell_{\text{3DGS}} + \ell_{\text{image}}]. \quad (7)$$

While the 3D supervision signal from the noisy teacher is not perfect, it is already in the 3D domain, making it computationally efficient. This setup allows for standard single-step denoising training, which is faster and less memory-intensive, with additional robustness introduced by the image-based supervision. Training the diffusion model in this way brings it to a performance level comparable to the base teacher model, preparing it for the multi-step training stage, where it can be fine-tuned to significantly surpass the base model's performance.

3.3. Cycle Consistency Regularization

Both the bootstrapping and fine-tuning phases with multi-step denoising utilize the image rendering loss. Inspired by cycle consistency losses in unpaired image-to-image translation [67], we propose to further regularize the model using the generated output \hat{s}_0 by utilizing the rendered image $\hat{x}_{\text{tgt}} = \mathcal{R}(\hat{s}_0, v_{\text{tgt}})$ to drive a second Gaussian Splats prediction, denoted as \tilde{s}_0 . We then render this second prediction back to the source view to define our cycle consistency loss term:

$$\mathcal{L}_{\text{cyc}} = \|x_{\text{src}} - \mathcal{R}(\tilde{s}_0, v_{\text{src}})\|_2^2. \quad (8)$$

Intuitively, this loss aims to constrain the predicted rendered view not only to match the target image in terms of appearance similarity, but also to be reliable enough to drive the generation of the source view. This loss is applied in both training stages. In the bootstrapping phase, the second splat prediction \tilde{s}_0 is generated through the noisy teacher, maintaining efficiency by only requiring one additional network pass. As shown in our ablation study, this loss improves the performance of the bootstrapping phase. We note that this technique could, in principle, also be used to improve the base model used as the noisy teacher, although this is beyond the scope of this work.

In the multi-step fine-tuning phase, however, our model already outperforms the noisy teacher (even without the cycle consistency loss), so lifting the predicted image to 3D via the noisy teacher is not meaningful. Instead, we apply the multi-step denoising process directly.

4. Experiments

4.1. Experimental Setups

Datasets. We conduct experiments using two datasets: the object-level ShapeNet-SRN [6, 44] and the scene-level RealEstate10k [66]. ShapeNet-SRN comprises synthetic objects across various categories. In line with Splatter Image [48] and PixelNeRF [61], we focus on the *cars* and *chairs* classes. The resolution for ShapeNet-SRN dataset is 128×128 , and the Splatter Image model is employed as the teacher for the ShapeNet experiments. RealEstate10k consists of real-world video data captured in both indoor and outdoor environments. Following Flash3D [47], we use a resolution of 256×384 for training in our experiments. The Flash3D model serves as the teacher to guide our diffusion model at the bootstrapping stage.

Evaluation Metrics. We adopt PSNR, SSIM [52] and LPIPS [64] as metrics for the evaluation of the image reconstruction and novel view synthesis.

Memory Usage and Model Size We report both GPU memory consumption and model size. As shown in Tab.2, our model exhibits a significantly smaller size compared to VisionNeRF and Splatter Image. While PixelNeRF has a

| Method | 1-view Cars | | | 1-view Chairs | | |
|------------------------------------|-----------------|-----------------|--------------------|-----------------|-----------------|--------------------|
| | PSNR \uparrow | SSIM \uparrow | LPIPS \downarrow | PSNR \uparrow | SSIM \uparrow | LPIPS \downarrow |
| SRN [44] | 22.25 | 0.88 | 0.129 | 22.89 | 0.89 | 0.104 |
| CodeNeRF [18] | 23.80 | 0.91 | 0.128 | 23.66 | 0.90 | 0.166 |
| FE-NVS [13] | 22.83 | 0.91 | 0.099 | 23.21 | 0.92 | 0.077 |
| ViewsetDiff w/o \mathcal{D} [46] | 23.21 | 0.90 | 0.116 | 24.16 | 0.91 | 0.088 |
| ViewsetDiff w \mathcal{D} [46] | 23.29 | 0.91 | 0.094 | - | - | - |
| PixelNeRF [61] | 23.17 | 0.89 | 0.146 | 23.72 | 0.90 | 0.128 |
| VisionNeRF [27] | 22.88 | 0.90 | 0.084 | 24.48 | 0.92 | 0.077 |
| NeRFDiff [12] | 23.95 | 0.92 | 0.092 | 24.80 | 0.93 | 0.070 |
| Splatter Image [48] | 24.00 | 0.92 | 0.078 | 24.43 | 0.93 | 0.067 |
| SplatDiffusion | 24.84 | 0.93 | 0.077 | 25.21 | 0.93 | 0.066 |

Table 1. **ShapeNet-SRN: Single-View Reconstruction (test split)**. Our method achieves better quality on all metrics on the Car split and Chair dataset, while performing reconstruction in the 3D space.

| Method | Memory Usage (GB) | Model Size (MB) |
|---------------------|-------------------|-----------------|
| PixelNeRF [61] | 3.05 | 113 |
| VisionNeRF [27] | 6.42 | 1390 |
| Splatter Image [48] | 1.71 | 646 |
| Ours | 1.15 | 295 |

Table 2. Memory Footprint and Model Size.

smaller model size, our approach achieves lower GPU memory consumption on the ShapeNet-SRN dataset.

4.2. Implementation Details.

Multi-step Denoising. We train the model using 4 NVIDIA A6000 GPUs. The computational efficiency is demonstrated in Tab. 2. During the bootstrapping stage (stage 1), a batch size of 100 per GPU is employed to train the diffusion model under the guidance of the teacher model. Following this, in stage 2, multi-step denoising is performed using a DDIM sampler with 10 inference steps. To manage the increased computational complexity during this phase, the batch size is reduced to 10.

Misc. We use Adam [23] as our optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$. We use a total noising steps of 100, with a linear scheduling, starting from 0.0001 to 0.2. For the bootstrapping stage, we use the teacher model to provide both supervision and noised samples. Instead of predicting the noises added to the splatters, our diffusion model denoises the noised inputs to clean samples directly. We set t^* to be 20. For the architecture of diffusion model, we use the U-Net implementation from diffusers¹. For the consistency branch, we use a denoising step of 10.

¹<https://huggingface.co/docs/diffusers>

4.3. Image conditioned reconstruction

ShapeNet-SRN. We benchmark our diffusion model on the ShapeNet-SRN dataset, as presented in Tab. 1. Using only a single input view, our model achieves PSNR improvements of 0.84 and 0.78 on the cars and chairs splits, respectively, compared to the Splatter Image baseline.

For qualitative evaluation, we compare our method with Splatter Image, which serves as our teacher model in Fig. 1 and in Fig. 3. As seen in the first row (Fig. 3 (a)), images generated by Splatter Image occasionally exhibit artifacts and distortions. In contrast, our model generally produces more fine-grained structures and higher-quality details. Furthermore, as shown in Fig. reffig:qual (b), the Gaussians generated by our model are denser and exhibit regular shapes, whereas those produced by Splatter Image tend to be oversized and less uniform.

RealEstate10K. We evaluate our method against recent state-of-the-art approaches on the real-world RealEstate10K dataset. As shown in Tab. 3, our model outperforms the teacher network, Flash3D, across three different evaluation settings, achieving an average PSNR improvement of 0.5. The visual comparisons in Fig. 1 and Fig. 3 further demonstrates the superiority of our method, consistently producing cleaner images while Flash3D struggles in unseen regions, resulting in blurry artifacts.

4.4. Additional View Guidance

Unlike deterministic feedforward models, diffusion models have the distinct advantage of incorporating guidance. In our approach, we condition the prediction of Gaussian Splat parameters on a single input view and can optionally leverage a second view as guidance during the denoising process, following the Universal Guidance framework [3]. Detailed explanations and formulations of the guidance mechanism

| Model | 5 frames | | | 10 frames | | | $\mathcal{U}[-30, 30]$ frames | | |
|---------------------|-----------------|-----------------|--------------------|-----------------|-----------------|--------------------|-------------------------------|-----------------|--------------------|
| | PSNR \uparrow | SSIM \uparrow | LPIPS \downarrow | PSNR \uparrow | SSIM \uparrow | LPIPS \downarrow | PSNR \uparrow | SSIM \uparrow | LPIPS \downarrow |
| Syn-Sin [55] | - | - | - | - | - | - | 22.30 | 0.740 | - |
| SV-MPI [50] | 27.10 | 0.870 | - | 24.40 | 0.812 | - | 23.52 | 0.785 | - |
| BTS [56] | - | - | - | - | - | - | 24.00 | 0.755 | 0.194 |
| Splatter Image [48] | 28.15 | 0.894 | 0.110 | 25.34 | 0.842 | 0.144 | 24.15 | 0.810 | 0.177 |
| MINE [26] | 28.45 | 0.897 | 0.111 | 25.89 | 0.850 | 0.150 | 24.75 | 0.820 | 0.179 |
| Flash3D [47] | 28.46 | 0.899 | 0.100 | 25.94 | 0.857 | 0.133 | 24.93 | 0.833 | 0.160 |
| SplatDiffusion | 29.12 | 0.932 | 0.087 | 26.54 | 0.887 | 0.122 | 25.40 | 0.873 | 0.135 |

Table 3. **Novel View Synthesis.** Our model shows superior performance on RealEstate10k on small, medium and large baseline ranges.

| Setting | Novel view synthesis | | | Source view synthesis | | |
|---|----------------------|-----------------|--------------------|-----------------------|-----------------|--------------------|
| | PSNR \uparrow | SSIM \uparrow | LPIPS \downarrow | PSNR \uparrow | SSIM \uparrow | LPIPS \downarrow |
| (a.1) Splatter Image | 24.1992 | 0.9213 | 0.0843 | 31.1158 | 0.9808 | 0.0269 |
| (a.2) feedforwad (our architecture) | 19.9947 | 0.8613 | 0.1588 | 23.2363 | 0.9165 | 0.0955 |
| (b.1) stage I w only rendering loss | 18.8201 | 0.8415 | 0.1862 | 20.9767 | 0.8815 | 0.1535 |
| (b.2) stage I w diffusion & rendering loss | 22.6078 | 0.9046 | 0.1083 | 28.2025 | 0.9690 | 0.0411 |
| (b.3) stage II w diffusion & rendering loss | 23.1323 | 0.9116 | 0.1061 | 29.4463 | 0.9750 | 0.0358 |
| (b.4) stage II w only rendering loss | 24.4936 | 0.9264 | 0.0945 | 31.9839 | 0.9850 | 0.0233 |
| (c.1) stage II w timestep from $[0, T]$ | 24.6998 | 0.9233 | 0.0907 | 33.1196 | 0.9820 | 0.0205 |
| (c.2) stage II w timestep from $[t^*, T]$ | 24.8324 | 0.9300 | 0.0871 | 33.4658 | 0.9883 | 0.0176 |
| (d.1) stage I w/o consistency | 22.6078 | 0.9046 | 0.1083 | 28.2025 | 0.9690 | 0.0411 |
| (d.2) stage I w consistency | 23.7293 | 0.9181 | 0.0979 | 29.9227 | 0.9774 | 0.0254 |
| (d.3) stage I, stage II w consistency | 24.9137 | 0.9332 | 0.0847 | 33.7061 | 0.9886 | 0.0153 |

Table 4. Ablations Studies on Single view Reconstruction, evaluated on the validation set of ShapeNet-SRN Cars

are provided in the supplementary material.

Table 5 compares our view guidance method to a 2-view 3DGS optimization procedure, as outlined by [21], which is initialized using the base model. Our diffusion model demonstrates a 0.2 PSNR improvement when incorporating image guidance, with an additional 0.2 PSNR gain achieved through Gaussian Splits optimization, consistently outperforming the Splatter Image baseline. While here we demonstrate guidance in a two-view settings, the guidance mechanism can naturally be extended to multiview scenarios, which we leave for future exploration.

4.5. Ablation

We conducted a series of ablation studies on the ShapeNet-SRN cars dataset to measure the effect of various architectural designs on both novel and source view synthesis. The results are summarized in Tab. 4.

Architectural Comparison. To assess whether our improvements stem from the diffusion framework or architectural changes, we trained a feedforward model using the same U-Net architecture as our diffusion model, with the timestep-

| Method | GS optim | Guidance | PSNR | SSIM | LPIPS |
|----------------|--------------|--------------|-------|------|-------|
| Splatter Image | \times | \times | 24.75 | 0.93 | 0.06 |
| Ours | \checkmark | \times | 25.24 | 0.94 | 0.06 |
| Ours | \times | \times | 25.18 | 0.93 | 0.06 |
| Ours | \times | \checkmark | 25.36 | 0.94 | 0.06 |
| Ours | \checkmark | \checkmark | 25.55 | 0.95 | 0.05 |

Table 5. **Additional-view guidance.** Evaluated on a subset of the car split, our diffusion-based model better utilizes an additional view through guidance compared to 3DGS optimization.

related layers removed to create a deterministic version. This feedforward network directly predict gaussian parameter from the input image as Splatter Image. Due to a much smaller model size, it performs significantly worse than Splatter Image (Tab. 4.a), underscoring that the diffusion approach is essential to the enhanced results we observe.

Bootstrapping (stage 1). Bootstrapping is necessary for the initialization of our diffusion model. As shown in Tab. 4(b.1), it produces unsatisfactory results to directly train diffusion

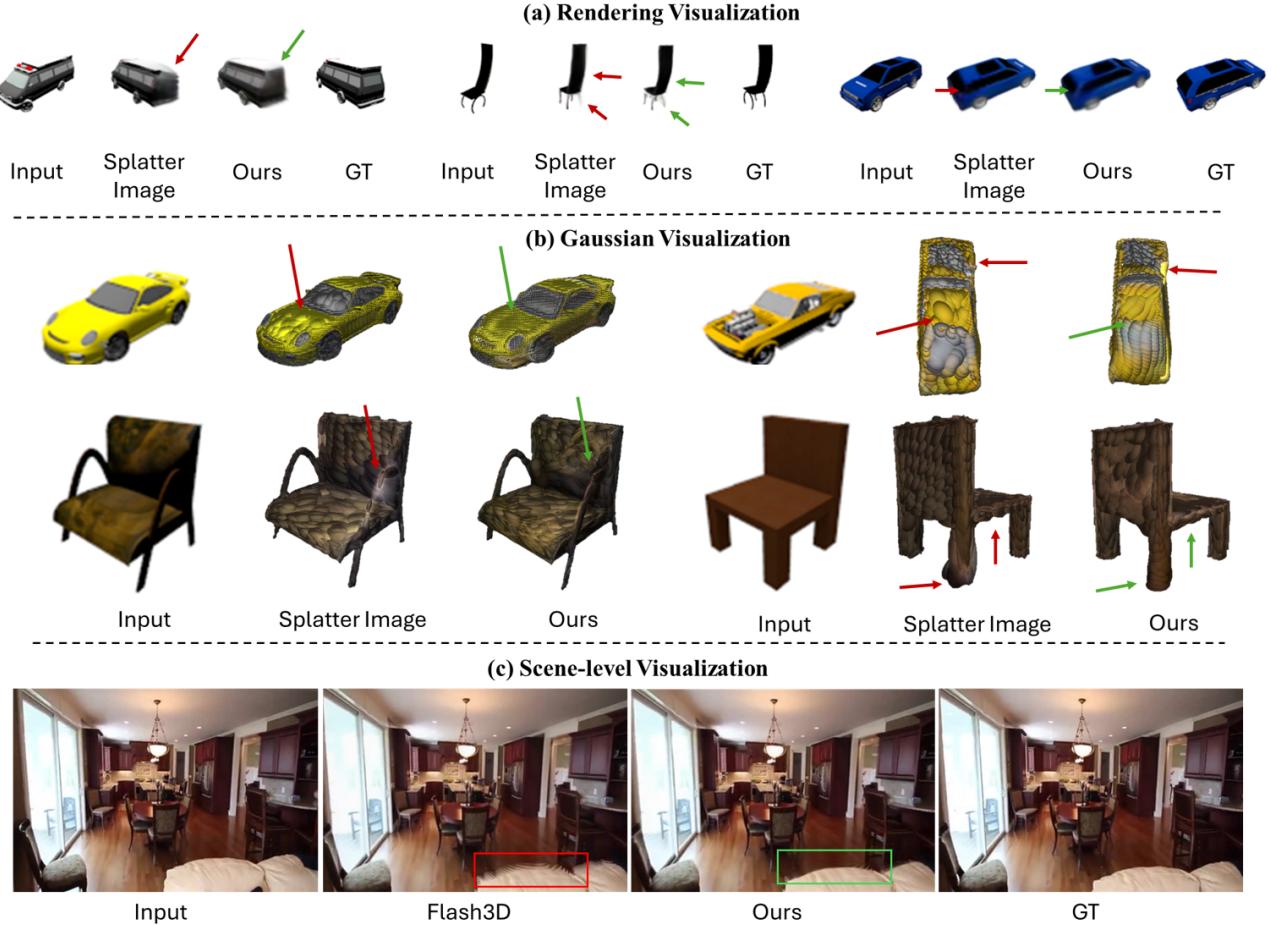


Figure 3. **Qualitative results.** (a) Qualitative comparison on the ShapeNet-SRN dataset for the Car and Chair categories. Our model produces views that are more faithful to the source image and better maintain plausibility. (b) Comparison of Gaussian Splat outputs between Splatter Image and our diffusion model shows that our model generates more regular patterns that closely follow the object surface. (c) Scene-level qualitative comparison on the RealEstate10K dataset demonstrates that our method produces more realistic results, particularly in ambiguous areas, such as the 2D edge separating the bed and the floor.

model without the teacher model as guidance because of the indirect cross-modality supervision. With the teacher guidance, the diffusion model can produce better results (Tab.4 (b.2)), but still bounded by the teacher’s performance.

Multi-step denoising (stage 2). In Stage 2, we found that the teacher model limits the performance of our model if we continue to use it as guidance (Tab.4 (b.3)) Instead, we fine-tune the model only with the rendering loss, allowing the model to explore how to improve the rendering performance from the ground truth images.

Cycle consistency. By introducing a feedback loop in which the predicted target view images are rendered back to the source view and supervised with the ground truth input image, we achieve performance improvements in both stages, as demonstrated in Tab. 4(d).

5. Conclusion and Limitations

In this work, we introduced a novel framework for training 3D diffusion models without requiring large-scale 3D datasets. By leveraging deterministic predictors as noisy teachers and using sparse 2D views for supervision, our approach enables effective training of 3D diffusion models with significant performance improvements.

Limitations. Our framework is flexible and could extend to various 3D representations; however, the current implementation relies on pixel-aligned 3D GS, inheriting certain limitations. Specifically, the uneven Gaussian distribution—where Gaussians concentrate on visible views with insufficient coverage in occluded regions—can lead to oversmoothness in novel views. Future work could address this limitation by adapting our framework to support alternative 3D representations, further enhancing its robustness and generalizability.

Acknowledgment

Or Litany is a Taub fellow and is supported by the Azrieli Foundation Early Career Faculty Fellowship. This research was supported in part by an academic gift from Meta.

References

- [1] Antonio Alliegro, Yawar Siddiqui, Tatiana Tommasi, and Matthias Nießner. Polydiff: Generating 3d polygonal meshes with diffusion models. *arXiv preprint arXiv:2312.11417*, 2023. [2](#) [3](#)
- [2] Titas Auciukevičius, Zexiang Xu, Matthew Fisher, Paul Henderson, Hakan Bilen, Niloy J Mitra, and Paul Guerrero. Renderdiffusion: Image diffusion for 3d reconstruction, inpainting and generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12608–12618, 2023. [4](#)
- [3] Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Universal guidance for diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 843–852, 2023. [6](#) [3](#)
- [4] Miguel Angel Bautista, Pengsheng Guo, Samira Abnar, Walter Talbott, Alexander Toshev, Zhuoyuan Chen, Laurent Dinh, Shuangfei Zhai, Hanlin Goh, Daniel Ulbricht, et al. Gaudi: A neural architect for immersive 3d scene generation. *Advances in Neural Information Processing Systems*, 35:25102–25116, 2022. [3](#)
- [5] Eric R. Chan, Koki Nagano, Matthew A. Chan, Alexander W. Bergman, Jeong Joon Park, Axel Levy, Miika Aittala, Shalini De Mello, Tero Karras, and Gordon Wetzstein. GeNVS: Generative novel view synthesis with 3D-aware diffusion models. In *arXiv*, 2023. [4](#)
- [6] 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. [5](#)
- [7] Hansheng Chen, Jiatao Gu, Anpei Chen, Wei Tian, Zhuowen Tu, Lingjie Liu, and Hao Su. Single-stage diffusion nerf: A unified approach to 3d generation and reconstruction. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2416–2425, 2023. [3](#)
- [8] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. [2](#)
- [9] Emilien Dupont, Hyunjik Kim, SM Eslami, Danilo Rezende, and Dan Rosenbaum. From data to functa: Your data point is a function and you can treat it like one. *arXiv preprint arXiv:2201.12204*, 2022. [3](#)
- [10] Bradley Efron. Tweedie’s formula and selection bias. *Journal of the American Statistical Association*, 106(496):1602–1614, 2011. [2](#)
- [11] Ruiqi Gao*, Aleksander Holynski*, Philipp Henzler, Arthur Brussee, Ricardo Martin-Brualla, Pratul P. Srinivasan, Jonathan T. Barron, and Ben Poole*. Cat3d: Create anything in 3d with multi-view diffusion models. *Advances in Neural Information Processing Systems*, 2024. [4](#)
- [12] Jiatao Gu, Alex Trevithick, Kai-En Lin, Joshua M Susskind, Christian Theobalt, Lingjie Liu, and Ravi Ramamoorthi. Nerfdiff: Single-image view synthesis with nerf-guided distillation from 3d-aware diffusion. In *International Conference on Machine Learning*, pages 11808–11826. PMLR, 2023. [6](#)
- [13] Pengsheng Guo, Miguel Angel Bautista, Alex Colburn, Liang Yang, Daniel Ulbricht, Joshua M Susskind, and Qi Shan. Fast and explicit neural view synthesis. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3791–3800, 2022. [6](#)
- [14] Amir Hertz, Kfir Aberman, and Daniel Cohen-Or. Delta denoising score. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2328–2337, 2023. [3](#)
- [15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. [2](#)
- [16] Lukas Höller, Aljaž Božič, Norman Müller, David Novotny, Hung-Yu Tseng, Christian Richardt, Michael Zollhöfer, and Matthias Nießner. Viewdiff: 3d-consistent image generation with text-to-image models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5043–5052, 2024. [4](#)
- [17] Yicong Hong, Kai Zhang, Juxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli, Trung Bui, and Hao Tan. Lrm: Large reconstruction model for single image to 3d. *arXiv preprint arXiv:2311.04400*, 2023. [3](#)
- [18] Wonbong Jang and Lourdes Agapito. Codenerf: Disentangled neural radiance fields for object categories. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12949–12958, 2021. [6](#)
- [19] Animesh Karnewar, Andrea Vedaldi, David Novotny, and Niloy J Mitra. Holodiffusion: Training a 3d diffusion model using 2d images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18423–18433, 2023. [4](#)
- [20] Oren Katzir, Or Patashnik, Daniel Cohen-Or, and Dani Lischinski. Noise-free score distillation. *arXiv preprint arXiv:2310.17590*, 2023. [3](#)
- [21] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. [7](#)
- [22] Jaihoon Kim, Juil Koo, Kyeongmin Yeo, and Minhyuk Sung. Sync tweedies: A general generative framework based on synchronized diffusions. *arXiv preprint arXiv:2403.14370*, 2024. [2](#)
- [23] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. [6](#)
- [24] Jeong-gi Kwak, Erqun Dong, Yuhe Jin, Hanseok Ko, Shweta Mahajan, and Kwang Moo Yi. Vivid-1-to-3: Novel view synthesis with video diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6775–6785, 2024. [4](#)
- [25] Kyungmin Lee, Kihyuk Sohn, and Jinwoo Shin. Dreamflow: High-quality text-to-3d generation by approximating probability flow. *arXiv preprint arXiv:2403.14966*, 2024. [3](#)

- [26] Jiaxin Li, Zijian Feng, Qi She, Henghui Ding, Changhu Wang, and Gim Hee Lee. Mine: Towards continuous depth mpi with nerf for novel view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12578–12588, 2021. 7, 2
- [27] Kai-En Lin, Yen-Chen Lin, Wei-Sheng Lai, Tsung-Yi Lin, Yi-Chang Shih, and Ravi Ramamoorthi. Vision transformer for nerf-based view synthesis from a single input image. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 806–815, 2023. 6, 1
- [28] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9298–9309, 2023. 4
- [29] Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. Syncdreamer: Generating multiview-consistent images from a single-view image. *arXiv preprint arXiv:2309.03453*, 2023. 4
- [30] Zhen Liu, Yao Feng, Michael J Black, Derek Nowrouzezahrai, Liam Paull, and Weiyang Liu. Meshdiffusion: Score-based generative 3d mesh modeling. *arXiv preprint arXiv:2303.08133*, 2023. 2, 3
- [31] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2837–2845, 2021. 3
- [32] David McAllister, Songwei Ge, Jia-Bin Huang, David W. Jacobs, Alexei A. Efros, Aleksander Holynski, and Angjoo Kanazawa. Rethinking score distillation as a bridge between image distributions. *arXiv preprint arXiv:2406.09417*, 2024. 3
- [33] 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
- [34] Yuxuan Mu, Xinxin Zuo, Chuan Guo, Yilin Wang, Juwei Lu, Xiaofeng Wu, Songcen Xu, Peng Dai, Youliang Yan, and Li Cheng. Gsd: View-guided gaussian splatting diffusion for 3d reconstruction. *arXiv preprint arXiv:2407.04237*, 2024. 2, 3
- [35] Norman Müller, Yawar Siddiqui, Lorenzo Porzi, Samuel Rota Bulò, Peter Kontschieder, and Matthias Nießner. Diffrrf: Rendering-guided 3d radiance field diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4328–4338, 2023. 3
- [36] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 3
- [37] Guocheng Qian, Jinjie Mai, Abdullah Hamdi, Jian Ren, Aliaksandr Siarohin, Bing Li, Hsin-Ying Lee, Ivan Skorokhodov, Peter Wonka, Sergey Tulyakov, and Bernard Ghanem. Magic123: One image to high-quality 3d object generation using both 2d and 3d diffusion priors. In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024. 3
- [38] Davis Rempe, Zhengyi Luo, Xue Bin Peng, Ye Yuan, Kris Kitani, Karsten Kreis, Sanja Fidler, and Or Litany. Trace and pace: Controllable pedestrian animation via guided trajectory diffusion. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2
- [39] Barbara Roessle, Norman Müller, Lorenzo Porzi, Samuel Rota Bulò, Peter Kotschieder, Angela Dai, and Matthias Nießner. L3dg: Latent 3d gaussian diffusion. *arXiv preprint arXiv:2410.13530*, 2024. 2, 3
- [40] 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
- [41] 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. 2
- [42] Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. Mvdream: Multi-view diffusion for 3d generation, 2024. 4
- [43] J Ryan Shue, Eric Ryan Chan, Ryan Po, Zachary Ankner, Jiajun Wu, and Gordon Wetzstein. 3d neural field generation using triplane diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20875–20886, 2023. 3
- [44] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. *Advances in Neural Information Processing Systems*, 32, 2019. 5, 6, 2
- [45] Ido Sobol, Chenfeng Xu, and Or Litany. Zero-to-hero: Enhancing zero-shot novel view synthesis via attention map filtering. *arXiv preprint arXiv:2405.18677*, 2024. 4
- [46] Stanislaw Szymanowicz, Christian Rupprecht, and Andrea Vedaldi. Viewset diffusion:(0-) image-conditioned 3d generative models from 2d data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8863–8873, 2023. 2, 4, 6
- [47] Stanislaw Szymanowicz, Eldar Insafutdinov, Chuanxia Zheng, Dylan Campbell, João F Henriques, Christian Rupprecht, and Andrea Vedaldi. Flash3d: Feed-forward generalisable 3d scene reconstruction from a single image. *arXiv preprint arXiv:2406.04343*, 2024. 3, 5, 7
- [48] Stanislaw Szymanowicz, Christian Rupprecht, and Andrea Vedaldi. Splatter image: Ultra-fast single-view 3d reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10208–10217, 2024. 3, 5, 6, 7
- [49] Ayush Tewari, Tianwei Yin, George Cazenavette, Semon Rezhnikov, Josh Tenenbaum, Frédéric Durand, Bill Freeman, and Vincent Sitzmann. Diffusion with forward models: Solving stochastic inverse problems without direct supervision. *Advances in Neural Information Processing Systems*, 36: 12349–12362, 2023. 4
- [50] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 551–560, 2020. 7

- [51] Arash Vahdat, Francis Williams, Zan Gojcic, Or Litany, Sanja Fidler, Karsten Kreis, et al. Lion: Latent point diffusion models for 3d shape generation. *Advances in Neural Information Processing Systems*, 35:10021–10039, 2022. 2, 3
- [52] Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4): 600–612, 2004. 5
- [53] 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. *Advances in Neural Information Processing Systems*, 36, 2024. 3
- [54] Daniel Watson, William Chan, Ricardo Martin Bruealla, Jonathan Ho, Andrea Tagliasacchi, and Mohammad Norouzi. Novel view synthesis with diffusion models. In *The Eleventh International Conference on Learning Representations*, 2023. 4
- [55] Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. Synsin: End-to-end view synthesis from a single image. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7467–7477, 2020. 7
- [56] Felix Wimbauer, Nan Yang, Christian Rupprecht, and Daniel Cremers. Behind the scenes: Density fields for single view reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9076–9086, 2023. 7
- [57] Rundi Wu, Ben Mildenhall, Philipp Henzler, Keunhong Park, Ruiqi Gao, Daniel Watson, Pratul P. Srinivasan, Dor Verbin, Jonathan T. Barron, Ben Poole, and Aleksander Holynski. Reconfusion: 3d reconstruction with diffusion priors. *arXiv*, 2023. 4
- [58] Chenfeng Xu, Huan Ling, Sanja Fidler, and Or Litany. 3diff-tion: 3d object detection with geometry-aware diffusion features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10617–10627, 2024.
- [59] Jiayu Yang, Ziang Cheng, Yunfei Duan, Pan Ji, and Hong-dong Li. Consistnet: Enforcing 3d consistency for multi-view images diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7079–7088, 2024. 4
- [60] Mingyang Yi, Aoxue Li, Yi Xin, and Zhenguo Li. Towards understanding the working mechanism of text-to-image diffusion model. *arXiv preprint arXiv:2405.15330*, 2024. 2, 4
- [61] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4578–4587, 2021. 5, 6, 1
- [62] 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. 3
- [63] Kai Zhang, Sai Bi, Hao Tan, Yuanbo Xiangli, Nanxuan Zhao, Kalyan Sunkavalli, and Zexiang Xu. Gs-lrm: Large reconstruction model for 3d gaussian splatting. *European Conference on Computer Vision*, 2024. 3
- [64] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric, 2018. 5
- [65] Linqi Zhou, Yilun Du, and Jiajun Wu. 3d shape generation and completion through point-voxel diffusion. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5826–5835, 2021. 3
- [66] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817*, 2018. 5
- [67] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017. 5
- [68] Zi-Xin Zou, Zhipeng Yu, Yuan-Chen Guo, Yangguang Li, Ding Liang, Yan-Pei Cao, and Song-Hai Zhang. Triplane meets gaussian splatting: Fast and generalizable single-view 3d reconstruction with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10324–10335, 2024. 2

A Lesson in Splats: Teacher-Guided Diffusion for 3D Gaussian Splats Generation with 2D Supervision

Supplementary Material

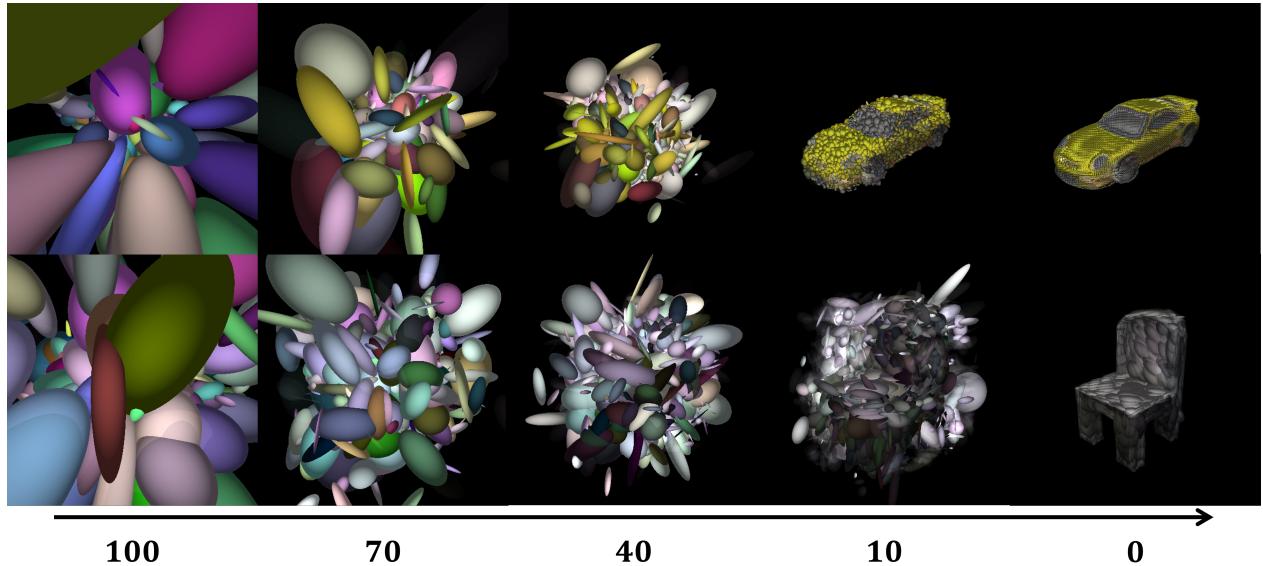


Figure 4. Visualization of the denoising process of our diffusion models, trained on the Car and Chair categories of ShapeNet-SRN dataset.

6. Additional results

6.1. Additional qualitative results

We present visual comparisons of our method to Pixel-NeRF [61] and VisionNeRF [27] on ShapeNet-SRN Cars and Chairs in Fig. 5. More qualitative results from RealEstate10K dataset is in Fig. 6.

6.2. Additional ablations

Feedforward vs Diffusion Model. To evaluate whether the observed improvements originated from the diffusion framework or architectural modifications, we trained a feed-forward model by removing the time-conditioning layers from the U-Net architecture while preserving its overall structure. For comparison, we predicted Gaussian parameters from a single input image following the splatter image. The feedforward model exhibited significantly worse performance, which we attribute to its reduced size, resulting in limited representational capacity. From Tab. 6, we conclude that the diffusion framework is more suitable for such generation tasks compared to deterministic models, producing better results even with a smaller model size.

Choices of losses. Through experiments (Tab. 7), we found that it produces terrible results to directly train a diffusion model using rendering loss (both in stage 1 and stage 2), because the supervision indirectly comes from the rendered

| Setting | PSNR \uparrow | SSIM \uparrow | LPIPS \downarrow |
|----------------|-----------------|-----------------|--------------------|
| Splatter Image | 24.1992 | 0.9213 | 0.0843 |
| Feedforwad | 19.9947 | 0.8613 | 0.1588 |

Table 6. Feedforward model vs Splatter Image.

image instead of the denoised splatters, which makes it hard for the diffusion model to learn the accurate distribution.

| Stage 1 | Stage 2 | PSNR \uparrow | SSIM \uparrow | LPIPS \downarrow |
|----------|----------|-----------------|-----------------|--------------------|
| \times | R | 16.7284 | 0.7836 | 0.3733 |
| R | \times | 18.8201 | 0.8415 | 0.1862 |
| D | \times | 21.3050 | 0.8965 | 0.1182 |
| R + D | \times | 22.6078 | 0.9046 | 0.1083 |
| R + D | R + D | 23.1323 | 0.9116 | 0.1061 |
| R + D | R | 24.4936 | 0.9264 | 0.0945 |

Table 7. Ablation of losses at two stages. ‘R’ and ‘D’ represent rendering loss and diffusion loss, respectively.

For stage 1 training, the performance improves using teacher model as guidance and it reports the best results using both rendering loss and diffusion loss.

For stage 2 training, if we continue to use the diffusion

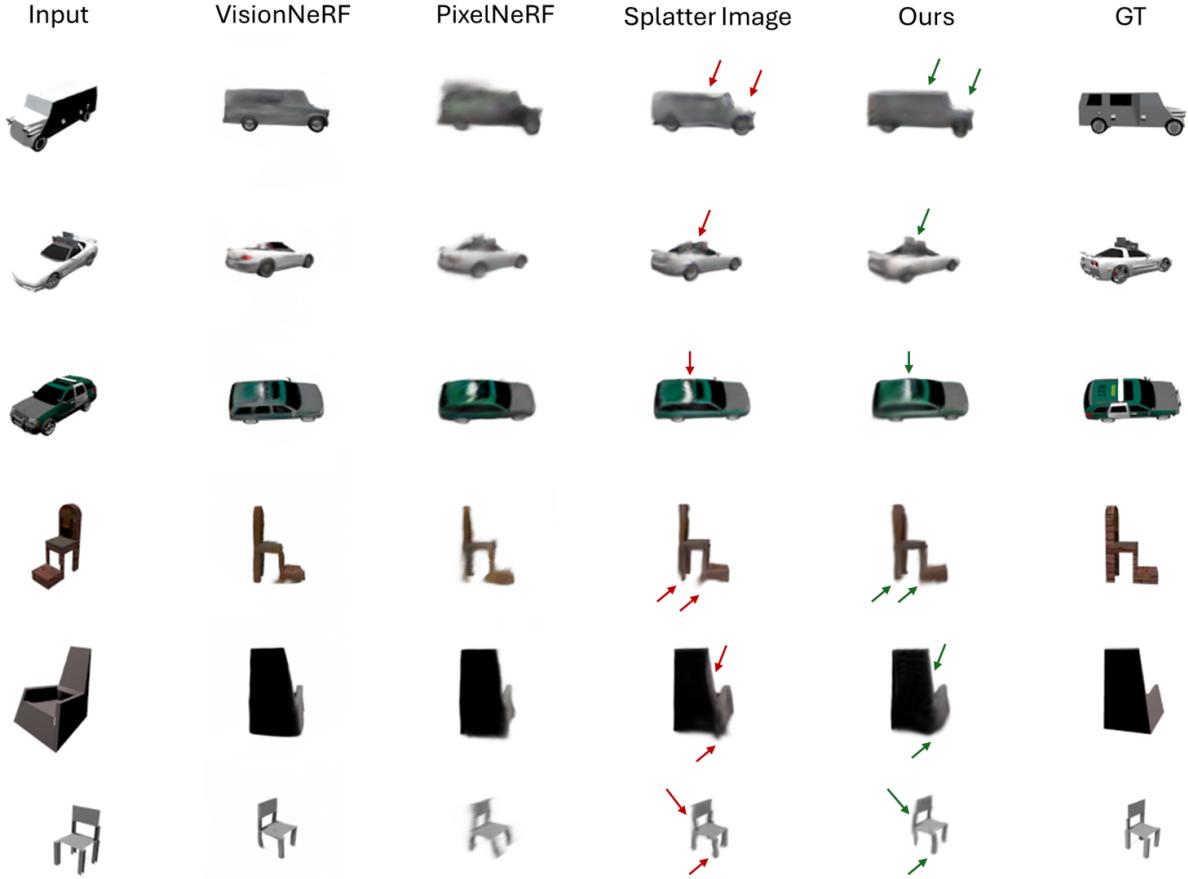


Figure 5. **Additional qualitative results.** Qualitative comparisons on the ShapeNet-SRN dataset for additional viewpoints and objects from the Car and Chair categories. Our model produces views that are more faithful to the source image and better maintain plausibility, while maintaining the fast rendering of Splatter Image.

loss, the teacher model will limit the performance of our diffusion model. Therefore, we only use rendering loss at stage 2, allowing the model to explore how to minimize the rendering loss and improve the rendering performance.

Weighted loss at different timesteps. The difficulty of prediction at different timesteps varies. Therefore, during the stage 2 training, we assign different weights to the rendering loss obtained at different timesteps and accumulate them for back-propagation throughout the denoising steps. The ablation results are in Tab. 8.

| Setting | PSNR \uparrow | SSIM \uparrow | LPIPS \downarrow |
|-------------------|-----------------|-----------------|--------------------|
| w/o weighted loss | 22.8848 | 0.9116 | 0.1044 |
| w weighted loss | 24.4936 | 0.9264 | 0.0945 |

Table 8. Ablations of weighted loss at different timesteps

7. Data details

7.1. ShapeNet-SRN Cars and Chairs

We adhere to the standard protocol for the ShapeNet-SRN dataset. Specifically, we use the provided images, camera intrinsics, camera poses, and data splits provided by [44] with a resolution of 128×128 . Our method is trained using relative camera poses. For single-view reconstruction, view 64 serves as the conditioning view, while for additional-view guidance, view 88 is used as guidance view. All remaining available views are treated as target views, where we compute novel view synthesis metrics.

7.2. RealEstate10K

We obtain 65,384 videos and their corresponding camera pose trajectories from the provided youtube links. Using these camera poses, we perform sparse point cloud reconstruction with COLMAP [41]. For evaluation, we adopt the test split provided by MINE [26] and follow prior work by



Figure 6. **Additional qualitative results.** Qualitative comparisons on RealEstate10K dataset.

assessing PSNR on novel frames that are 5 and 10 frames ahead of the source frame. Additionally, we evaluate on a randomly sampled frame within an interval of ± 30 frames, using the same frames employed in MINE’s evaluation. For evaluation, we use a total of 3,205 frames. The results presented in Tab. 3 are sourced from Flash3D [47]. Our model is trained and tested at a resolution of 256×384 .

8. Implementation details

Multi-step Denoising. We train the model on 4 NVIDIA A6000 GPUs. Our diffusion model is quite efficient . For bootstrapping at stage 1, we use a batch size of 100 on each GPU. After obtaining the diffusion model from the teacher model, we perform multi-step denoising with a DDIM sampler of 10 inference steps. The batch size for stage 2 reduces to 10. We assign different weights to the rendering loss obtained at different timesteps and accumulate them for back-propagation throughout the denoising steps.

Additional-view guidance Different from deterministic feedforward models, one significant advantage we gain from diffusion models is the ability of using guidance. We use one input view as the condition to predict the Gaussian Splats parameters and then use a second view as guidance during the denoising process using the forward guidance from Universal Guidance [3].

Since we predict \hat{s}_0 directly, the noise can be calculated as follows:

$$\epsilon_t = \frac{s_t - \sqrt{\alpha_t} \hat{s}_0}{\sqrt{1 - \alpha_t}} \quad (9)$$

Then we calculate the gradient using the guidance image x_{gd} and the corresponding view direction v :

$$\text{grad} \leftarrow \nabla_{s_t} \ell[x_{gd}, \mathcal{R}(\hat{s}_0, v)]. \quad (10)$$

With the guidance strength factor $s(t)$, we can obtain $\hat{\epsilon}_t$

$$\hat{\epsilon}_t = \epsilon_t + s(t) \cdot \text{grad} \quad (11)$$

At last, we can get s_{t-1} following DDIM sampling:

$$s_{t-1} = \sqrt{\alpha_{t-1}} \hat{s}_0 + \sqrt{1 - \alpha_{t-1}} \cdot \hat{\epsilon}_t \quad (12)$$

| Method | GS optim | Guidance | PSNR | SSIM | LPIPS |
|----------|----------|----------|-------|------|-------|
| Splatter | ✗ | ✗ | 24.75 | 0.93 | 0.06 |
| Image | ✓ | ✗ | 25.24 | 0.94 | 0.06 |
| | ✗ | ✗ | 25.18 | 0.93 | 0.06 |
| Ours | ✓ | ✗ | 25.26 | 0.94 | 0.06 |
| | ✗ | ✓ | 25.36 | 0.94 | 0.06 |
| | ✓ | ✓ | 25.55 | 0.95 | 0.05 |

Table 9. **Additional-view guidance.** Evaluated on a subset of the car split, because per-sample GS optimization takes time.