
DiffGS: Functional Gaussian Splatting Diffusion

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 3D Gaussian Splatting (3DGS) has shown convincing performance in rendering
2 speed and fidelity, yet the generation of Gaussian Splatting remains a challenge due
3 to its discreteness and unstructural nature. In this work, we propose DiffGS, a gen-
4 eral Gaussian generator based on latent diffusion models. DiffGS is a powerful and
5 efficient 3D generative model which is capable of generating Gaussian primitives
6 at arbitrary numbers for high-fidelity rendering with rasterization. The key insight
7 is to disentangledly represent Gaussian Splatting via three novel functions to model
8 Gaussian probabilities, colors and transforms. Through the novel disentangling of
9 3DGS, we represent the discrete and unstructural 3DGS with continuous Gaussian
10 Splatting functions, where we then train a latent diffusion model with the target
11 of generating these Gaussian Splatting functions both unconditionally and condi-
12 tionally. Meanwhile, we introduce a discretization algorithm to extract Gaussians
13 at arbitrary numbers from the generated functions via octree-guided sampling
14 and optimization. We explore DiffGS for various tasks, including unconditional
15 generation, conditional generation from text, image, and partial 3DGS, as well as
16 Point-to-Gaussian generation. We believe that DiffGS provides a new direction for
17 flexibly modeling and generating Gaussian Splatting.

18

1 Introduction

19 3D content creation is a vital task in computer graphics and 3D computer vision, which shows
20 great potential in real-world applications such as virtual reality, game design, film production, and
21 robotics. Previous 3D generative models usually take Neural Radiance Field (NeRF) [30, 2, 47] as
22 the representation. However, the volumetric rendering for NeRF requires considerable computational
23 cost, leading to sluggish rendering speeds and significant memory burden. Recent advances of 3D
24 Gaussian Splatting (3DGS) [21, 52, 19] demonstrating its potential to serve as the next-generation 3D
25 representation by enabling both real-time rendering and high-fidelity appearance modeling. Designing
26 3D generative models for 3DGS provides a scheme for real-time interaction with 3D creations.

27 The core challenge in generative 3DGS modeling lies in its discreteness and unstructural nature, which
28 prevents the well-studied frameworks in structural image/voxel/video generation from transferring to
29 directly generate 3DGS. Concurrent works [56, 15] alternatively transport Gaussians into structural
30 voxel grids with volume generation models [9] for generating Gaussians. However, these methods
31 lead to 1) abundant computational cost for high-resolution voxels, 2) struggles in preserving high-
32 quality Gaussian reconstructions due to the information loss during voxelization and 3) limited
33 number of generated Gaussians constrained by the voxel resolutions.

34 To address these challenges, we present DiffGS, a novel diffusion-based generative model for general
35 3D Gaussian Splatting which is capable of efficiently generating high-quality Gaussian primitives
36 at arbitrary numbers. The key insight of DiffGS is to disentangledly represent Gaussian Splatting
37 via three novel functions: Gaussian probability function (GauPF), Gaussian color function (GauCF)
38 and Gaussian transform function (GauTF). GauPF indicates the geometry of 3DGS by modeling the

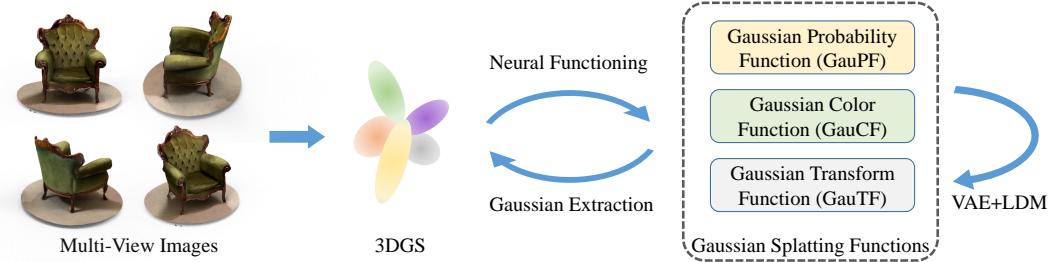


Figure 1: Illustration of DiffGS. We fit 3DGS from multi-view images and then disentangle it into three Gaussian Splating Functions. We train a Gaussian VAE with a latent diffusion model for generating these functions, followed by a Gaussian extraction algorithm to obtain the final generated Gaussians.

39 probabilities of each sampled 3D location to be a Gaussian location. GauCF and GauTF predict the
 40 Gaussian attributes of appearances and transformations given a 3D location as input. Through the
 41 novel disentangling of 3DGS, we represent the discrete and unstructural 3DGS with three continuous
 42 Gaussian Splating Functions.

43 With the disentangled and powerful representation, the next step is to design a generative model
 44 with the target of generating these Gaussian Splating functions. We propose a Gaussian VAE model
 45 for creating compressed representation for the Gaussian Splating functions. The Gaussian VAE
 46 learns a regularized latent space which maps the Gaussians Splating functions of each shape into
 47 one latent vector. A latent diffusion model (LDM) is simultaneously trained at the latent space for
 48 generating novel 3DGS shapes. With the powerful LDM, We explore DiffGS to generate diverse
 49 3DGS both conditionally and unconditionally. Finally, we introduce a discretization algorithm to
 50 extract Gaussians at arbitrary numbers from the generated functions via octree-guided sampling and
 51 optimization. The key idea is to first extract 3D Gaussian geometry from GauPF as the 3D locations
 52 sampled at the 3D spaces with largest Gaussian probabilities, and then predict the Gaussian attributes
 53 with GauCF and GauTF. We illustrate the overview of DiffGS in Fig. 1.

54 We systematically summarize the superiority of DiffGS in terms of: 1) Efficiency, we design DiffGS
 55 based on Gaussian Splating and Latent Diffusion Models, which shows significant efficiency in
 56 model training, inference and shape rendering. 2),3) Generality and quality, we generate native 3DGS
 57 without processes like voxelization, leading to unimpaired quality and generality in applying to
 58 downstream 3DGS applications. 4) Scalability, we scalably generate Gaussian primitives at arbitrary
 59 numbers. We conduct comprehensive experiments on both synthetic ShapeNet dataset and real-world
 60 DeepFashion3D dataset, which demonstrate our non-trivial improvements over the state-of-the-art
 61 methods.

62 In summary, our contributions are:

- 63 • We propose DiffGS, a novel diffusion-based generative model for general 3D Gaussian
 64 Splating which is capable of efficiently generating high-quality Gaussian primitives at
 65 arbitrary numbers.
- 66 • We introduce a novel schema to disentangledly represent Gaussian Splating via three
 67 functions to model Gaussian probabilities, Gaussian colors and Gaussian transforms. We
 68 simultaneously propose a discretization algorithm to extract Gaussians from these functions
 69 via octree-guided sampling and optimization.
- 70 • DiffGS achieves remarkable performances under various tasks including unconditional
 71 generation, conditional generation from text, image, and partial 3DGS, as well as Point-to-
 72 Gaussian generation.

73 2 Related Work

74 2.1 Rendering-Guided 3D Representation

75 Remarkable progress have been achieved in the field of novel view synthesis (NVS) [30, 34, 47, 2,
 76 32], since the proposal of Neural Radiance Field (NeRF) [30]. NeRF implicitly represents scene

77 appearance and geometries using MLP-based neural networks, optimized through volume rendering
78 to achieve outstanding novel view synthesis quality. Some variants [1, 12, 36] show promising
79 performance by advancing NeRF in terms of rendering quality, scalability and view-consistency.
80 Additionally, more recent methods [32, 6, 11, 49] explore the training and rendering efficiency of
81 NeRF by introducing feature-grids based 3D representations. Instant-NGP [32] highly accelerates
82 NeRF learning by introducing multiresolution feature grids based on hash table with fully-fused
83 CUDA kernel implementations. However, the NeRF representations which require expensive neural
84 network inferences during volume rendering, still struggles in the applications where real-time
85 rendering speed is required.

86 Recently, the emergence of 3D Gaussian Splatting (3DGS) [21, 44, 23, 52, 55] has showcased
87 impressive real-time results in novel view synthesis (NVS). 3DGS [21] has led to revolutions in the
88 NVS field by demonstrating superior performances in multiple domains. However, the generation of
89 Gaussian Splatting remains a challenge due to its discreteness and unstructural nature. We introduce
90 a novel schema to represent the discrete and unstructural 3DGS with three continuous Gaussian
91 Splatting Functions, thus ingeniously tackle the challenge by designing generative models for the
92 functions.

93 2.2 3D Generative Models

94 The field of creating 3D contents with generative models has emerged as a particularly captivating
95 research direction. A series of works [35, 24, 50, 39, 29, 27, 7, 43, 53, 45] focus on optimization-
96 based frameworks based on Score Distillation Sampling (SDS), which achieve convincing generation
97 performances by distilling 3D geometry and appearance of the radiance fields with pretrained 2D
98 diffusion models [33, 17] as the prior. However, these works entail significant computational costs
99 due to time-consuming per-scene optimization. Going beyond optimization-based 3D generation,
100 recent works [31, 46, 48, 20] explore 3D generative methods based on diffusion models to directly
101 learn priors from 3D datasets for generative radiance fields modeling. These methods typically
102 represent radiance fields as structural triplanes [48, 42, 14] or voxels [46, 31, 9]. DiffRF [31] leverage
103 a voxel based NeRF representation with 3D U-Nets as the backbone to train a diffusion model.

104 With the recent advances in 3DGS [21], designing a powerful 3D generative model for generating
105 3DGS is expected to be a popular research topic and also brings significant challenges due to
106 the discreteness and unstructural nature of 3DGS, which prevents the well-studied frameworks
107 in structural image/voxel/video generation from transferring to directly generate 3DGS. A series
108 of works [54, 59, 18, 57] follow the schema of image-based reconstruction without generative
109 modeling, which lack the ability to generate diverse shapes. Concurrent works [15, 56] follow
110 the voxel-based representations to transport Gaussians into structural voxel grids with volume
111 generation models for generating Gaussians. However, these methods come with several drawbacks,
112 including abundant computational costs for high-resolution voxels, struggles in preserving high-
113 quality Gaussian reconstructions due to information loss during voxelization, and a restricted number
114 of generated Gaussians constrained by voxel resolutions. In contrast, our proposed DiffGS explore
115 a new perspective to directly represent the discrete and unstructural 3DGS with three continuous
116 Gaussian Splatting Functions. Though the insight, we design a latent diffusion model for efficiently
117 generating high-quality Gaussian primitives by learning to generate the Gaussian Splatting Functions.
118 DiffGS generates general Gaussians at arbitrary numbers with a specially designed octree-based
119 extraction algorithm.

120 3 Method

121 We introduce DiffGS, a novel diffusion-based generative model for general 3D Gaussian Splatting
122 which is capable of efficiently generating high-quality Gaussian primitives at arbitrary numbers. The
123 overview of DiffGS is shown in Fig. 2. We first preview Gaussian Splatting in Sec. 3.1 and present
124 the novel functional schema for representing Gaussian Splatting with three disentangled Gaussian
125 Splatting Functions in Sec. 3.2. We then introduce the Gaussian Variational Auto-encoder and the
126 Latent Diffusion model for compressing and generative modeling on Gaussian Splatting Functions,
127 as shown in Sec. 3.3. A novel discretization algorithm is proposed in Sec. 3.4 to extract Gaussians at
128 arbitrary numbers from the generated functions via octree-guided sampling and optimization.

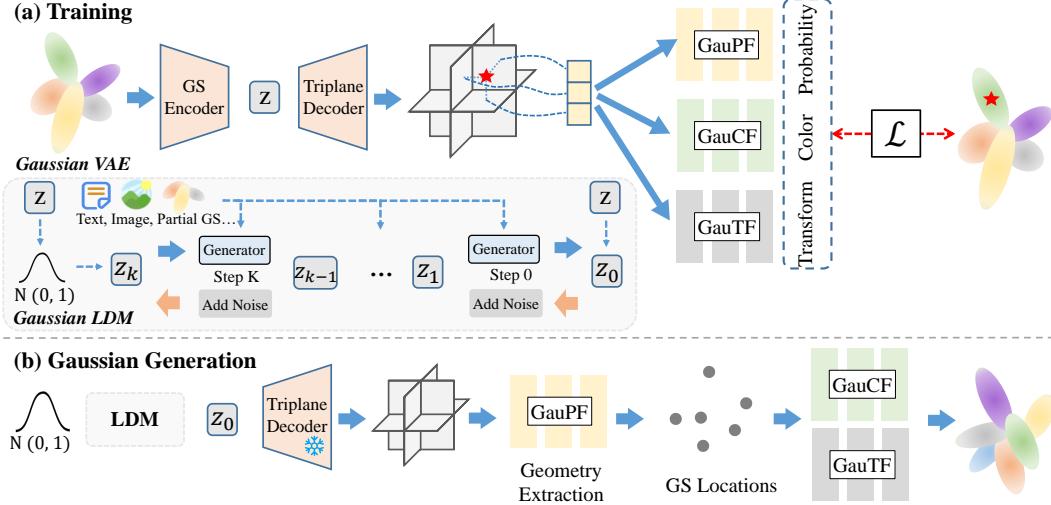


Figure 2: Overview of DiffGS. (a) We disentangle the fitted 3DGS into three Gaussian Splatting Functions to model the Gaussian probability, colors and transforms. We then train a Gaussian VAE with a conditional latent diffusion model for generating these functions. (b) During generation, we first extract Gaussian geometry from the generated GauPF, followed by the GauCF and GauTF to obtain the Gaussian attributes.

129 3.1 Preview Gaussian Splatting

130 3D Gaussian Splatting (3DGS) [21] represents a 3D shape or scene as a set of Gaussians with
 131 attributes to model the geometries and view-dependent appearances. For a 3DGS $G = \{g_i\}_{j=1}^N$
 132 containing N Gaussians, the geometry of j -th Gaussian is explicitly parameterized via 3D covariance
 133 matrix Σ_i and its center $\sigma_i \in \mathbb{R}^3$, formulated as:

$$g_i(x) = \exp\left(-\frac{1}{2}(x - \sigma_i)^T \Sigma^{-1} (x - \sigma_i)\right), \quad (1)$$

134 where the covariance matrix $\Sigma_i = r_i s_i s_i^T r_i^T$ is factorized into a rotation matrix $r_i \in \mathbb{R}^4$ and a scale
 135 matrix $s_i \in \mathbb{R}$. The appearance of the Gaussian g_i is controlled by an opacity value $o_j \in \mathbb{R}$ and a
 136 color value $c_i \in \mathbb{R}^3$. Note that the color is represented as a series of sphere harmonics coefficients
 137 in practice of 3DGS, yet we still keep its definition as three-dimension color c_i in our paper for a
 138 clear understanding on our method. To this end, the 3DGS G is defined as $\{g_i = \{\sigma_i, r_i, s_i, o_i, c_i\} \in \mathbb{R}^K\}_{j=1}^N$, where K is dimension of the combined attributes in each Gaussian.
 139

140 3.2 Functional Gaussian Splatting Representation

141 The key challenge in generative 3DGS modeling lies in its discreteness and unstructural nature, which
 142 prevents the well-studied generative frameworks from transferring to directly generate 3DGS. We
 143 address this challenge by introducing to disentangledly represent Gaussian Splatting via three novel
 144 functions: Gaussian probability function (GauPF), Gaussian color function (GauCF) and Gaussian
 145 transform function (GauTF). Through the novel disentangling of 3DGS, we represent the discrete
 146 and unstructural 3DGS with three continuous Gaussian Splatting Functions.

147 **Gaussian Probability Function.** Gaussian Probability Function (GauPF) indicates the geometry
 148 of 3DGS by modeling the probabilities of each sampled 3D location to be a Gaussian location.
 149 Given a set of 3D query location $Q = \{q_j \in \mathbb{R}^3\}_{j=1}^M$ sampled in 3D space around a fitted 3DGS
 150 $G = \{g_i \in \mathbb{R}^3\}_{j=1}^N$, the GauPF of G predicts the probabilities p of queries $\{q_j\}_{j=1}^M$ to be a Gaussian
 151 location in G , formulated as:

$$p_j = \text{GauPF}(q_j) \in [0, 1]. \quad (2)$$

152 The idea of Gaussian probability modeling comes from the observation that the further a 3D location
 153 q_j is from all Gaussians, the lower the probability that any Gaussian occupies the space at q_j .

154 Therefore the ground truth Gaussian probability of q_j is defined as:

$$\text{GauPF}(q_j) = \tau(\lambda(\min_{i \in [1, N]} \|\hat{q}_j - \sigma_i\|_2)), \quad (3)$$

155 where $\min_{i \in [1, N]} \|\hat{q}_j - \sigma_i\|_2$ indicates the distance from q_j to the nearest Gaussian center in $\{\sigma_i\}_{i=1}^N$, λ
 156 is a truncation function which filters the extremely large values and τ is a continuous function which
 157 maps the query-to-Gaussian distances to probabilities in the range of [0,1].

158 A learned GauPF implicitly models the locations of 3D Gaussian centers, which is the key factor for
 159 generating high-quality 3DGS. The extraction of 3DGS centers from GauPF is then achieved with
 160 our designed Gaussian extraction algorithm which will be introduced in Sec. 3.4.

161 **Gaussian Color and Transform Modeling.** Gaussian Color Function (GauCF) and Gaussian
 162 Transform Function (GauTF) predict the Gaussian attributes of appearances and transformations from
 163 Gaussian geometries. Specifically, given the center σ_i of a Gaussian g_i in G as input, GauCF predicts
 164 the color attribute c_i and GauTF predicts the rotation r_i , scale s_i and opacity o_i , formulated as:

$$\{c_i\} = \text{GauCF}(\sigma_i); \quad \{r_i, s_i, o_i\} = \text{GauTF}(\sigma_i). \quad (4)$$

165 Note that GauCF and GauTF mainly focus on predicting the Gaussian colors and transforms from 3D
 166 Gaussian centers. This is different from the GauPF which models the probabilities of query samples
 167 in the 3D space. The reason is that GauPF focuses on exploring the geometry of 3DGS from the 3D
 168 space, while GauCF and GauTF learn to predict the Gaussian attributes from the known geometries.
 169 Through the novel disentangling of 3DGS, we represent the discrete and unstructural 3DGS with three
 170 continuous Gaussian Splatting Functions. The functional representation is a general and flexible term
 171 for 3DGS which has no restrictions on the Gaussian numbers, densities, geometries, etc.

172 3.3 Gaussian Variational Auto-encoder and Latent Diffusion

173 With the disentangled and powerful representation, the next step is to design a generative model
 174 with the target of generating these Gaussian Splatting functions. We follow the common schema to
 175 design a Gaussian Variational Auto-encoder (VAE) [22] with a Latent Diffusion Model (LDM) [33]
 176 as the generative model. The detailed framework and the training pipeline of DiffGS are illustrated in
 177 Fig. 1(a).

178 **Gaussian VAE.** The Gaussian VAE compresses the Gaussian Splatting Functions into a regularized
 179 latent space by mapping the Gaussian Splatting Functions of each 3DGS shape into a latent vector,
 180 from which we can also recover the Gaussian Splatting Functions. Specifically, the Gaussian VAE
 181 consists of 1) a GS encoder ϕ_{en} to learn representations from 3DGS and encodes each 3DGS into a
 182 latent vector z , 2) a triplane decoder ϕ_{de} which decodes the latent z into a feature triplane, and 3)
 183 three neural predictors ψ_{pf} , ψ_{cf} and ψ_{tf} which serve as the implementation of GauPF, GauCF and
 184 GauTF to predict Gaussian probabilities, colors and transforms.

185 Given a fitted 3DGS $G = \{g_i\}_{i=1}^N$ as input, the GS encoder ϕ_{en} extracts a global latent feature z from
 186 G , which is then decoded into a feature triplane $t \in \mathbb{R}^{H \times W \times C \times 3}$ with the decoder ϕ_{de} , formulated
 187 as:

$$z = \phi_{en}(G); \quad t = \phi_{de}(z). \quad (5)$$

188 The triplane t consists of three orthogonal feature planes $\{t_{XY}, t_{XZ}, t_{YZ}\}$ which are aligned to the
 189 axex. For a 3D location q_j , we obtain its corresponding feature $f_j = \text{interp}(t, q_j)$ from the triplane
 190 t by projecting q_j onto the orthogonal feature planes and concatenating the tri-linear interpolated
 191 features at the three planes. We then predict the Gaussian probability, color and transform of q_j with
 192 the neural Gaussian Function predictors as:

$$\{\hat{p}_i\} = \psi_{pf}(f_j); \quad \{\hat{c}_i\} = \psi_{cf}(f_j); \quad \{\hat{r}_i, \hat{s}_i, \hat{o}_i\} = \psi_{tf}(f_j). \quad (6)$$

193 The Gaussian VAE is trained with the target of accurately predicting Gaussian attributes and robustly
 194 regularizing the latent space. In practice, the training objective is formulated as:

$$\mathcal{L}_{\text{VAE}} = \|\{\hat{p}, \hat{c}, \hat{r}, \hat{s}, \hat{o}\} - \{p, c, r, s, o\}\|_1 + \beta(D_{KL}(\mathcal{Q}_\phi(z|G)\|\mathcal{P}(z))). \quad (7)$$

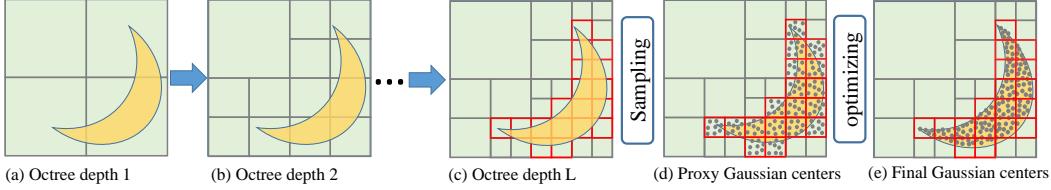


Figure 3: Gaussian geometry extractions from generated GauPF. The yellow and green regions indicate the high probability area and the low probability area judged by GauPF. (a),(b) and (c) show the progressively octree build process at depth 1,2 and L . (d) We sample proxy Gaussian centers from the octree at final depth L . (e) We optimize proxy centers to the exact geometry indicated in GauPF.

195 The first loss term indicates the \mathcal{L}_1 loss between the predicted Gaussian attributes in Eq. (6)
 196 and the target ones defined by the ground truth Gaussian Splatting Functions in Eq. (2) and Eq. (4).
 197 The second term in Eq. (7) is the KL-divergence loss with a factor of β , which constrains on the
 198 regularization of the learned latent space of z . Specifically, we define the inferred posterior of
 199 z as the distribution $Q_\phi(z|G)$, which is regularized to align with the Gaussian distribution prior
 200 $\mathcal{P}(z) = \mathcal{N}(0, I)$, where I is the standard deviation.

201 **Gaussian LDM.** With the trained Gaussian VAE in place, we are now able to encode any 3DGS
 202 into a compact 1D latent vector z . We then train a latent diffusion model (LDM) [33] efficiently
 203 on the latent space. A diffusion model is trained to generate samples from a target distribution by
 204 reversing a process that incrementally introduces noise. We define $\{z_0, z_1, \dots, z_K\}$ as the forward
 205 process $\gamma(z_{0:K})$ which gradually transforms a real data z_0 into Gaussian noise (z_T) by adding noises.
 206 The backward process $\mu(z_{0:K})$ leverages a neural generator μ to denoise z_K into a real data sample.

207 To achieve controlled generation of 3DGS, we introduce a conditioning mechanism [41] into the
 208 diffusion process with cross-attention. Given an input condition y (e.g. text, image, partial 3DGS),
 209 we leverage a custom encoder δ to project y into the condition embedding $\delta(y)$. The embedding is
 210 then fused into the generator μ with cross attention modules. Following DDPM, we simply adopt the
 211 optimizing objective to train the generator for predicting noises ϵ_σ , formulated as:

$$\mathcal{L}_{\text{LDM}} = \mathbb{E}_{z_0, t, \epsilon \sim \mathcal{N}(0, I)} \left[\|\epsilon - \epsilon_\sigma(z_t, \delta(y), t)\|^2 \right], \quad (8)$$

212 where t is a time step and ϵ is a noise latent sampled from the Gaussian distribution $\mathcal{N}(0, I)$. We
 213 adopt the well-studied architecture DALLE-2 [40] as the LDM implementation.

214 3.4 Gaussian Extraction Algorithm

215 The final step for the generation process of DiffGS is to extract 3DGS from the generated Gaussian
 216 Splatting Functions, similar to the effect of Marching Cubes algorithm [25] which extracts meshes
 217 from Signed Distance Functions. The key factor is to extract the geometries of 3DGS, i.e., Gaussian
 218 locations and the appearances of 3DGS, i.e., colors and transforms. The full generation pipeline is
 219 shown in Fig. 2(b).

220 **Octree-Guided Geometry Sampling.** The locations of 3D Gaussian centers indicate the geometry of
 221 the represented 3DGS. We aim to design a discretization algorithm to obtain the discrete 3D locations
 222 from the learned continuous Gaussian Probability Function parameterized with the neural network
 223 ψ_{pf} , which models the probability of each query sampled in the 3D space to be a 3D Gaussian
 224 location. To achieve this, we design an octree based sampling and optimization algorithm which
 225 generates accurate center locations of 3D Gaussians at arbitrary numbers.

226 We show the 2D illustration of the algorithm in Fig. 3. Assume that the 3D space is divided into
 227 the high probability area (the yellow region) and the low probability area (the green region) by
 228 the generated GauPF. We aim to extract the geometry as the locations with high probabilities. A
 229 naive implementation is to densely sample queries in the 3D space and keep the ones with large
 230 probabilities as outputs. However, it will lead to high computational cost for inferencing and the
 231 discrete sampling also struggles to accurately reach the locations with largest probabilities in the
 232 continuous GauPF. We get inspiration from octree [28, 51] to design a progressive strategy which
 233 only explores the 3D regions with large probabilities in current octree depth for further subdivision in
 234 the next octree depth. After L layers of octree subdivision, we reach the local regions with largest

Table 1: Comparisons of unconditional generation under ShapeNet [5] dataset.

Method	Airplane		Chair	
	FID-50K ↓	KID-50K (%) ↓	FID-50K ↓	KID-50K (%) ↓
GET3D [13]	—	—	59.51	2.414
DiffTF [4]	110.8	9.173	93.02	6.708
Ours	47.03	3.436	35.28	2.148



Figure 4: Visual comparisons with state-of-the-arts on unconditional generation of ShapeNet Chairs.

probabilities, from where we uniformly sample N 3D points as the proxy points $\{\rho_i\}_{i=1}^N$ representing coarse locations of Gaussian centers.

Optimizing Geometry with GauPF. To further refine the proxy points to the exact locations of Gaussian centers with largest probabilities in GauPF, we propose to further optimize the proxy points with the supervision from learned GauPF ψ_{pf} . Specifically, we set the position of proxy points $\{\rho_i = \{\rho_x, \rho_y, \rho_z\}\}_{i=1}^N$ to be learnable and optimize them to reach the positions $\{\hat{\sigma}_i\}_{i=1}^N$ with largest probabilities of ψ_{pf} . The optimization target is formulated as:

$$\mathcal{L}_{\text{Geo}} = \frac{1}{N} \sum_{i=1}^N \psi_{pf}(\rho_i). \quad (9)$$

Note that we can set N to arbitrary numbers, enabling DiffGS to generate 3DGS with no limits on the density and resolution.

Extracting Gaussian Attributes. We now obtain the estimated geometry indicating the predicted Gaussian centers $\{\hat{\sigma}_i\}_{i=1}^N$. We then extract the appearances and transforms from the generated triplane t , Gaussian Color Function ψ_{cf} and Gaussian Transform Function ψ_{tf} as $\{\hat{c}_i\} = \psi_{cf}(\text{interp}(t, \hat{\sigma}_i))$ and $\{\hat{r}_i, \hat{s}_i, \hat{o}_i\} = \psi_{tf}(\text{interp}(t, \hat{\sigma}_i))$. Finally, the general 3DGS is now generated as $\hat{G} = \{\hat{\sigma}_i, \hat{c}_i, \hat{r}_i, \hat{s}_i, \hat{o}_i\}_{i=1}^N$.

4 Experiment

4.1 Unconditional Generation

Dataset and Metrics. For unconditional generation of 3D Gaussian Splatting, we conduct experiments under the airplane and chair classes of ShapeNet [5] dataset. Following previous works [31, 4], we report two widely-used image generation metrics Fréchet Inception Distance (FID) [16] and Kernel Inception Distance (KID) [3] for evaluating the rendering quality of our proposed DiffGS and previous state-of-the-art works. The metrics are evaluated between 50K renderings of the generated shapes and 50K renderings of the ground truth ones, both at the resolution of 1024×1024 .

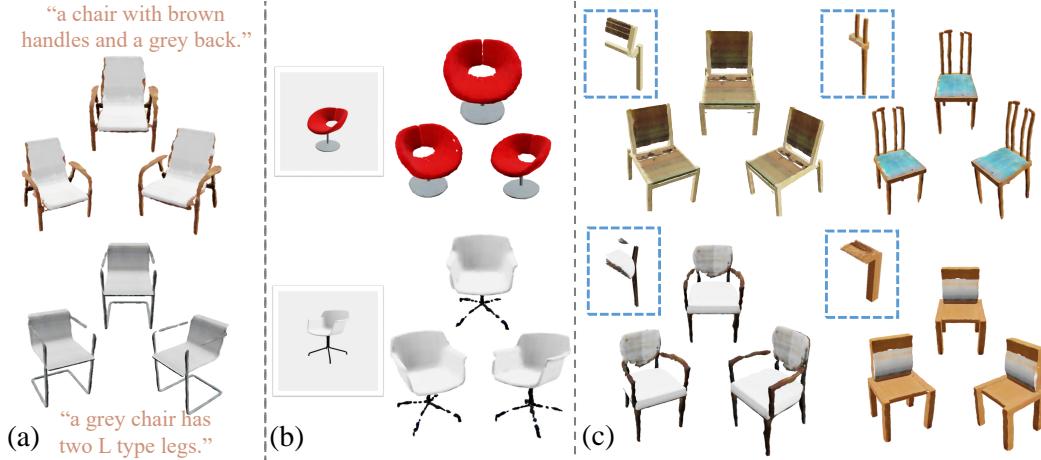


Figure 5: Visualization of conditional 3DGS generation results on ShapeNet. (a) Text conditional generation. (b) Image conditional generation. (c) Gaussian Splatting completion.

257 **Comparisons** We compare DiffGS with the state-of-the-art methods in terms of the rendering
 258 quality of generated shapes, including the GAN-based methods GET3D [13] and the diffusion based
 259 method DiffTF [4]. The quantitative comparison is shown in Tab. 1, where DiffGS achieves the best
 260 performance over all the baselines. We further show the visual comparison on the renderings of some
 261 generated shapes in Fig. 4, where the GAN-based GET3D struggle in generating complex shapes
 262 and the generations of DiffTF is blurry with poor textures. In contrast, DiffGS produces significantly
 263 more visual-appealing and high-fidelity generations in terms of rendering and geometry qualities.

264 4.2 Conditional Generation

265 We explore the conditional generation ability of DiffGS given texts, images and partial 3DGS as the
 266 input conditions. All the experiments are conducted under the chair class of ShapeNet [5] dataset
 267 with commonly used data splits in previous methods [8, 26].

268 **Text/Image-conditional Gaussian Splatting Generation** For introducing texts/images as the conditions
 269 for controllable Gaussian Splatting generation, we leverage the frozen text and image encoder
 270 from the pretrained CLIP [38] model as the implementation of custom text encoder γ_{text} and γ_{image}
 271 for achieving text/image embeddings. We then train DiffGS with the conditional optimization objective
 272 in Eq.(8). We show the visualization of some text/image conditional generations produced by
 273 DiffGS in Fig. 5(a) and Fig. 5(b). The results show that DiffGS accurately recovers the semantics
 274 and geometries described in the text prompts and the images, demonstrating the powerful capability
 275 of DiffGS in generating high-fidelity 3DGS from text descriptions or vision signals.

276 **Gaussian Splatting Completion** Additionally, we explore an interesting task of Gaussian Splatting
 277 completion. To the best of our knowledge, we are the first to focus and introduce solutions for this
 278 task. Specifically, the Gaussian Splatting completion task is to recover the complete 3DGS from a
 279 partial 3DGS which contains large occlusions. In real-world applications, having only sparse views
 280 with limited viewpoint movement available for optimizing 3DGS often results in a partial 3DGS.
 281 Solving Gaussian Splatting completion task enables us to infer the complete and dense 3DGS from
 282 the partial ones for improving the rendering quality at invisible viewpoints.

283 We introduce DiffGS with partial 3DGS as the conditions for solving this task. Specifically, we simply
 284 leverage a modified PointNet [37] as the custom encoder $\gamma_{partial}$ for partial 3DGS. Fig. 5(c) presents
 285 the visualization of Gaussian Splatting completion results produced by DiffGS. The results show that
 286 DiffGS is capable of recovering complex geometries and detailed appearances from highly occluded
 287 3DGS. Please refer to the appendix for implementation details on Gaussian Splatting completion.

288 4.3 Point-to-Gaussian Generation

289 We further introduce DiffGS for another challenging and vital task of Point-to-Gaussian generation.
 290 This task aims to generate the Gaussian attributes given a 3D point cloud as input. The task serves as

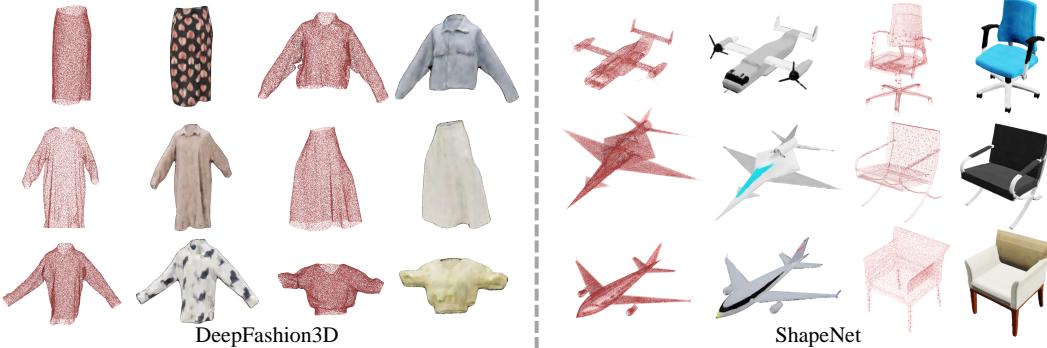


Figure 6: Visualization of Point-to-Gaussian fittings on DeepFashion3D and generations on ShapeNet.

the bridge between the easily accessible point clouds and the powerful 3DGS representation which efficiently models high-quality 3D appearances.

Dataset and Implementation. We conduct experiments under the chair and airplane classes of ShapeNet and also the widely-used garment dataset DeepFashion3D [58]. The DeepFashion3D dataset is a real-captured 3D dataset containing complex textures. For implementing Point-to-Gaussian, we simply train the Gaussian VAE with the three-dimension point clouds as inputs, instead of the 3DGS with attributes. Please refer to the appendix for more details on data preparation and implementation.

Performances. We provide the visualization of some Point-to-Gaussian fitting and generation results in Fig. 6. We show the fitting results for DeepFashion3D [58] dataset and the generation results for the test set of airplane and chair classes in ShapeNet [5]. DiffGS produces visual-appealing 3DGS generations given only 3D point cloud geometries as inputs. The results demonstrate that DiffGS can accurately predict Gaussian attributes for 3D point clouds. We believe DiffGS provides a new direction for 3DGS content generation by connecting 3DGS with point clouds.

4.4 Ablation Study

To evaluate some major designs and important hyper-parameters in DiffGS, we conduct ablation studies under the chair class of ShapeNet dataset. We report the performance in terms of PSNR, SSIM and LPIPS of the reconstructed 3DGS with Gaussian VAE.

Framework Designs. We first evaluate some major designs of our framework in Tab. 2. We justify the effectiveness of introducing the truncation function λ when modeling GauPF and report the results without λ as ‘w/o truncation’. We then explore implementing the projection function τ either as $\tau(x) = e^{-x}$ (as shown in ‘Exponent’) or as a linear projection (as shown in ‘Ours’). We also show the results without the optimization process during Gaussian extraction as ‘w/o Optimization’, which demonstrates the effectiveness of optimizing Gaussians to the exact locations.

Gaussian Numbers. One significant advantage of DiffGS lies in the ability of generating high-quality Gaussians at arbitrary numbers. To explore how the number of Gaussians affects the rendering quality, we conduct ablations on the Gaussian numbers as shown in Fig. 3. The results demonstrate that denser Gaussians lead to better quality.

Table 2: Ablations on framework designs.

Method	PSNR	SSIM	LPIPS
w/o Trunction	29.39	0.9792	0.0173
Exponent	29.74	0.9765	0.0188
w/o Optimization	30.34	0.9875	0.0152
Ours	34.01	0.9879	0.0149

Table 3: Ablations on Gaussian number.

Num	PSNR	SSIM	LPIPS
50K	28.61	0.9787	0.0251
100K	30.35	0.9838	0.0151
350K	34.01	0.9879	0.0149

5 Conclusion

In this paper, we introduce DiffGS for generative modeling of 3DGS. DiffGS disentangled represent 3DGS via three novel functions to model Gaussian probabilities, colors and transforms. We then train a latent diffusion model with the target of generating these functions both conditionally and unconditionally. DiffGS generates 3DGS with arbitrary numbers by an octree-guided extraction algorithm. The experimental results on various tasks demonstrate the superiority of DiffGS.

330

References

- [1] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5855–5864, 2021.
- [2] 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.
- [3] Mikołaj Bińkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. *arXiv preprint arXiv:1801.01401*, 2018.
- [4] Ziang Cao, Fangzhou Hong, Tong Wu, Liang Pan, and Ziwei Liu. Large-vocabulary 3d diffusion model with transformer. *arXiv preprint arXiv:2309.07920*, 2023.
- [5] 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.
- [6] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision*, pages 333–350. Springer, 2022.
- [7] Zilong Chen, Feng Wang, and Huaping Liu. Text-to-3d using gaussian splatting. *arXiv preprint arXiv:2309.16585*, 2023.
- [8] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5939–5948, 2019.
- [9] Yen-Chi Cheng, Hsin-Ying Lee, Sergey Tulyakov, Alexander G Schwing, and Liang-Yan Gui. Sdfusion: Multimodal 3d shape completion, reconstruction, and generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4456–4465, 2023.
- [10] Gene Chou, Yuval Bahat, and Felix Heide. Diffusion-sdf: Conditional generative modeling of signed distance functions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2262–2272, 2023.
- [11] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxtels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022.
- [12] Qiancheng Fu, Qingshan Xu, Yew-Soon Ong, and Wenbing Tao. Geo-Neus: Geometry-consistent neural implicit surfaces learning for multi-view reconstruction. *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- [13] Jun Gao, Tianchang Shen, Zian Wang, Wenzheng Chen, Kangxue Yin, Daiqing Li, Or Litany, Zan Gojcic, and Sanja Fidler. Get3d: A generative model of high quality 3d textured shapes learned from images. *Advances In Neural Information Processing Systems*, 35:31841–31854, 2022.
- [14] Anchit Gupta, Wenhan Xiong, Yixin Nie, Ian Jones, and Barlas Oğuz. 3dgen: Triplane latent diffusion for textured mesh generation. *arXiv preprint arXiv:2303.05371*, 2023.
- [15] Xianglong He, Junyi Chen, Sida Peng, Di Huang, Yangguang Li, Xiaoshui Huang, Chun Yuan, Wanli Ouyang, and Tong He. Gvgen: Text-to-3d generation with volumetric representation. *arXiv preprint arXiv:2403.12957*, 2024.
- [16] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- [17] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [18] Yicong Hong, Kai Zhang, Jiuxiang 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.
- [19] Binbin Huang, Zehao Yu, Anpei Chen, Andreas Geiger, and Shenghua Gao. 2d gaussian splatting for geometrically accurate radiance fields. *arXiv preprint arXiv:2403.17888*, 2024.
- [20] Heewoo Jun and Alex Nichol. Shap-e: Generating conditional 3d implicit functions. *arXiv preprint arXiv:2305.02463*, 2023.
- [21] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023.

- 389 [22] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- 390 [23] Jiahe Li, Jiawei Zhang, Xiao Bai, Jin Zheng, Xin Ning, Jun Zhou, and Lin Gu. Dngaussian:
392 Optimizing sparse-view 3d gaussian radiance fields with global-local depth normalization. *arXiv preprint arXiv:2403.06912*, 2024.
- 393 [24] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten
395 Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d
396 content creation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
397 Recognition*, pages 300–309, 2023.
- 398 [25] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3D surface
399 construction algorithm. *ACM Siggraph Computer Graphics*, 21(4):163–169, 1987.
- 400 [26] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In
401 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
402 2837–2845, 2021.
- 403 [27] Baorui Ma, Haoge Deng, Junsheng Zhou, Yu-Shen Liu, Tiejun Huang, and Xinlong Wang.
404 Geodream: Disentangling 2d and geometric priors for high-fidelity and consistent 3d generation.
405 *arXiv preprint arXiv:2311.17971*, 2023.
- 406 [28] Donald Meagher. Geometric modeling using octree encoding. *Computer graphics and image
407 processing*, 19(2):129–147, 1982.
- 408 [29] Gal Metzger, Elad Richardson, Or Patashnik, Raja Giryes, and Daniel Cohen-Or. Latent-nerf for
409 shape-guided generation of 3d shapes and textures. In *Proceedings of the IEEE/CVF Conference
410 on Computer Vision and Pattern Recognition*, pages 12663–12673, 2023.
- 411 [30] B Mildenhall, PP Srinivasan, M Tancik, JT Barron, R Ramamoorthi, and R Ng. NeRF:
412 Representing scenes as neural radiance fields for view synthesis. In *European Conference on
413 Computer Vision*, 2020.
- 414 [31] Norman Müller, Yawar Siddiqui, Lorenzo Porzi, Samuel Rota Bulo, Peter Kortschieder, and
415 Matthias Nießner. Diffrf: Rendering-guided 3d radiance field diffusion. In *Proceedings of the
416 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4328–4338, 2023.
- 417 [32] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics
418 primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–
419 15, 2022.
- 420 [33] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic
421 models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021.
- 422 [34] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M
423 Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *Proceedings
424 of the IEEE/CVF International Conference on Computer Vision*, pages 5865–5874, 2021.
- 425 [35] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using
426 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022.
- 427 [36] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf:
428 Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on
429 Computer Vision and Pattern Recognition*, pages 10318–10327, 2021.
- 430 [37] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point
431 sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer
432 vision and pattern recognition*, pages 652–660, 2017.
- 433 [38] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
434 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
435 models from natural language supervision. In *International conference on machine learning*,
436 pages 8748–8763. PMLR, 2021.
- 437 [39] Amit Raj, Srinivas Kaza, Ben Poole, Michael Niemeyer, Nataniel Ruiz, Ben Mildenhall, Shiran
438 Zada, Kfir Aberman, Michael Rubinstein, Jonathan Barron, et al. Dreambooth3d: Subject-
439 driven text-to-3d generation. In *Proceedings of the IEEE/CVF International Conference on
440 Computer Vision*, pages 2349–2359, 2023.
- 441 [40] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical
442 text-conditional image generation with clip latents, 2022.
- 443 [41] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
444 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF
445 conference on computer vision and pattern recognition*, pages 10684–10695, 2022.
- 446 [42] J Ryan Shue, Eric Ryan Chan, Ryan Po, Zachary Ankner, Jiajun Wu, and Gordon Wetzstein. 3d
447 neural field generation using triplane diffusion. In *Proceedings of the IEEE/CVF Conference on
448 Computer Vision and Pattern Recognition*, pages 20875–20886, 2023.

- 449 [43] Jingxiang Sun, Bo Zhang, Ruizhi Shao, Lizhen Wang, Wen Liu, Zhenda Xie, and Yebin Liu.
 450 Dreamcraft3d: Hierarchical 3d generation with bootstrapped diffusion prior. *arXiv preprint*
 451 *arXiv:2310.16818*, 2023.
- 452 [44] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative
 453 gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023.
- 454 [45] Junshu Tang, Tengfei Wang, Bo Zhang, Ting Zhang, Ran Yi, Lizhuang Ma, and Dong Chen.
 455 Make-it-3d: High-fidelity 3d creation from a single image with diffusion prior. In *Proceedings*
 456 *of the IEEE/CVF International Conference on Computer Vision*, pages 22819–22829, 2023.
- 457 [46] Zhicong Tang, Shuyang Gu, Chunyu Wang, Ting Zhang, Jianmin Bao, Dong Chen, and Baining
 458 Guo. Volumediffusion: Flexible text-to-3d generation with efficient volumetric encoder. *arXiv*
 459 *preprint arXiv:2312.11459*, 2023.
- 460 [47] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang.
 461 NeuS: Learning neural implicit surfaces by volume rendering for multi-view reconstruction.
 462 *Advances in Neural Information Processing Systems*, 34:27171–27183, 2021.
- 463 [48] Tengfei Wang, Bo Zhang, Ting Zhang, Shuyang Gu, Jianmin Bao, Tadas Baltrusaitis, Jingjing
 464 Shen, Dong Chen, Fang Wen, Qifeng Chen, et al. Rodin: A generative model for sculpting 3d
 465 digital avatars using diffusion. In *Proceedings of the IEEE/CVF conference on computer vision*
 466 *and pattern recognition*, pages 4563–4573, 2023.
- 467 [49] Yiming Wang, Qin Han, Marc Habermann, Kostas Daniilidis, Christian Theobalt, and Lingjie
 468 Liu. Neus2: Fast learning of neural implicit surfaces for multi-view reconstruction. In *Pro-*
 469 *ceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3295–3306,
 470 2023.
- 471 [50] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Pro-
 472 lificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation.
 473 *Advances in Neural Information Processing Systems*, 36, 2024.
- 474 [51] Jane Wilhelms and Allen Van Gelder. Octrees for faster isosurface generation. *ACM Transac-*
 475 *tions on Graphics (TOG)*, 11(3):201–227, 1992.
- 476 [52] Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi
 477 Tian, and Xinggang Wang. 4d gaussian splatting for real-time dynamic scene rendering. *arXiv*
 478 *preprint arXiv:2310.08528*, 2023.
- 479 [53] Jiale Xu, Xintao Wang, Weihao Cheng, Yan-Pei Cao, Ying Shan, Xiaohu Qie, and Shenghua
 480 Gao. Dream3d: Zero-shot text-to-3d synthesis using 3d shape prior and text-to-image diffu-
 481 sion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
 482 *Recognition*, pages 20908–20918, 2023.
- 483 [54] Yinghao Xu, Zifan Shi, Wang Yifan, Hansheng Chen, Ceyuan Yang, Sida Peng, Yujun Shen, and
 484 Gordon Wetzstein. Grm: Large gaussian reconstruction model for efficient 3d reconstruc-
 485 tion and generation. *arXiv preprint arXiv:2403.14621*, 2024.
- 486 [55] Taoran Yi, Jiemin Fang, Guanjun Wu, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Qi Tian, and
 487 Xinggang Wang. Gaussiandreamer: Fast generation from text to 3d gaussian splatting with
 488 point cloud priors. *arXiv preprint arXiv:2310.08529*, 2023.
- 489 [56] Bowen Zhang, Yiji Cheng, Jiaolong Yang, Chunyu Wang, Feng Zhao, Yansong Tang, Dong
 490 Chen, and Baining Guo. Gaussiancube: Structuring gaussian splatting using optimal transport
 491 for 3d generative modeling. *arXiv preprint arXiv:2403.19655*, 2024.
- 492 [57] Kai Zhang, Sai Bi, Hao Tan, Yuanbo Xiangli, Nanxuan Zhao, Kalyan Sunkavalli, and Zex-
 493 iang Xu. Gs-lrm: Large reconstruction model for 3d gaussian splatting. *arXiv preprint*
 494 *arXiv:2404.19702*, 2024.
- 495 [58] Heming Zhu, Yu Cao, Hang Jin, Weikai Chen, Dong Du, Zhangye Wang, Shuguang Cui, and
 496 Xiaoguang Han. Deep fashion3d: A dataset and benchmark for 3d garment reconstruction from
 497 single images. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK,*
 498 *August 23–28, 2020, Proceedings, Part I 16*, pages 512–530. Springer, 2020.
- 499 [59] Zi-Xin Zou, Zhipeng Yu, Yuan-Chen Guo, Yangguang Li, Ding Liang, Yan-Pei Cao, and
 500 Song-Hai Zhang. Triplane meets gaussian splatting: Fast and generalizable single-view 3d
 501 reconstruction with transformers. *arXiv preprint arXiv:2312.09147*, 2023.

502 **A Source Codes**

503 We provide our demonstration code as a part of our supplementary materials. We will release the
504 source code, data, pretrained models and instructions upon acceptance.

505 **B More Experimental Details**

506 **B.1 Gaussian Splatting Data Preparing**

507 DiffGS takes fitted 3DGS as input for learning generative modeling. To prepare the 3DGS dataset of
508 ShapeNet, we uniformly render 100 views from the ground truth meshes with blender to first obtain
509 the dense multi-view images for each 3D shape in the chair and airplane classes of ShapeNet dataset.
510 After that, we leverage the vanilla 3D Gaussian Splatting [21] method for fitting the 3DGS for each
511 shape with the rendered multi-view images.

512 For achieving more stable and regularized 3DGS data for better generative modeling, we design some
513 strategies for better initialization and optimization of 3DGS. (1) Since we fit 3DGS from existing
514 3D datasets with known geometries, we can simply sample dense point clouds uniformly from the
515 surfaces as the perfect initialization for 3DGS optimization, instead of initializing with COLMAP
516 points. The sampled point number is set to 100K. (2) We observe that optimizing 3DGS freely may
517 often lead to some extremely large Gaussians. This will lead to unstable training of the Gaussian VAE
518 and latent diffusion models, further affecting the generative modeling results. Therefore, we clip the
519 scales at a maximum size of 0.01 to avoid the abnormal Gaussians.

520 **B.2 Gaussian Splatting Completion**

521 We explore the task of Gaussian Splatting completion. Specifically, the Gaussian Splatting completion
522 task is to recover the complete 3DGS from a partial 3DGS which contains large occlusions. In
523 real-world applications, having only sparse views with limited viewpoint movement available for
524 optimizing 3DGS often results in a partial 3DGS. Solving Gaussian Splatting completion task enables
525 us to infer the complete and dense 3DGS from the partial ones for improving the rendering quality at
526 invisible viewpoints.

527 **Data preparation.** We generate partial 3D Gaussian Splatting data from the complete datasets in
528 a straightforward manner. First, we randomly divide each 3DGS into 8 chunks. Then, we occlude
529 7 chunks, leaving the remaining chunk as the partial 3DGS. This method allows us to prepare
530 partial-complete 3DGS pairs for training and testing.

531 **Implementation.** To leverage DiffGS with partial 3DGS as the conditions for Gaussian Splatting
532 completion, we leverage a modified PointNet [37] as the custom encoder $\gamma_{partial}$ for partial 3DGS,
533 which projects the partial 3DGS with K channels into a global partial 3DGS embedding. The DiffGS
534 for Gaussian Splatting completion is trained with the target of Eq.(8) by introducing partial 3DGS
535 embeddings through the cross-attention module.

536 **B.3 Point-to-Gaussian Generation**

537 **Data preparation.** For the task of Point-to-Gaussian generation, we first prepare the training/testing
538 data as point cloud-3DGS pairs obtained through the Gaussian Splatting Fitting process described in
539 Sec. B.1. Specifically, the point clouds are generated by densely sampling 100K points on the ground
540 truth meshes, while the paired 3DGS are obtained by optimizing with multi-view images rendered
541 around the ground truth meshes. The data preparation for the DeepFashion [58] dataset follows the
542 same process as for the ShapeNet [5] dataset.

543 Note that we train the Point-to-Gaussian DiffGS models for fitting the DeepFashion3D dataset, which
544 contains only 563 garment instances. In contrast, we split the airplane and chair classes of the
545 ShapeNet dataset into train/test sets to learn generalizable representations that enables DiffGS to
546 predict novel appearances for unknown point cloud geometries. The results shown in Fig. 6 of the
547 main paper include both the fitting results on the DeepFashion3D dataset and the generation results
548 from the point clouds in the test set of the airplane and chair classes in the ShapeNet dataset.

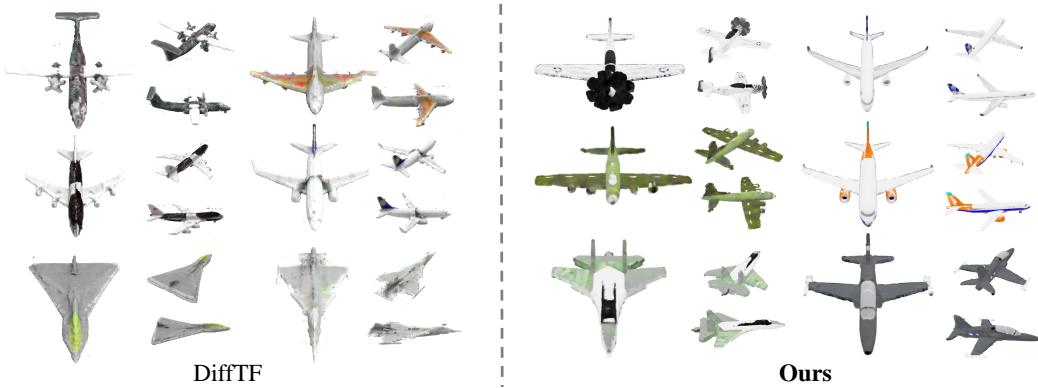


Figure 7: Visual comparisons with state-of-the-art method DiffTF [4] on unconditional generation of ShapeNet airplanes.

549 **Implementation.** For implementing Point-to-Gaussian, we train the Gaussian VAE using three-
 550 dimensional point clouds as inputs instead of 3DGS with attributes. Specifically, we replace the GS
 551 encoder in the Gaussian VAE with a PointNet-based network to learn representations from three-
 552 dimensional point clouds and recover Gaussian attributes from them. All architectures, optimization
 553 targets, and Gaussian extraction processes remain the same as in the Gaussian VAE, except for the
 554 encoder. Note that for the Point-to-Gaussian task, the Gaussian LDM is not trained, as we focus
 555 solely on decoding and extracting Gaussians from the point cloud inputs.

556 C Implement Details

557 We implement DiffGS with Pytorch Lightning. We leverage the Adam optimizer with a learning rate
 558 of 0.0001. We train DiffGS with eight 3090 GPUs and the convergence in each class of ShapeNet
 559 dataset takes around 5 days. The Guassian encoder and Triplane decoder are implemented based on
 560 the SDF-VAE encoder and decoder of DiffusionSDF [10], where we modify the PointNet [37] in the
 561 SDF-VAE encoder to receive K dimension 3DGS as the inputs.

562 D Additional Ablation Study on Octree Depth

563 The depth of the octree in our Gaussian extraction algo-
 564 rithm is an important hyper-parameter in the framework.
 565 To explore how the octree depth affects the rendering
 566 quality, we conduct ablations on the octree depth as
 567 shown in Fig. 4. The results demonstrate that a larger
 568 octree depth leads to better quality by capturing more
 569 geometry details.

Table 4: Ablations on Gaussian numbers.

Depth	PSNR	SSIM	LPIPS
7	29.77	0.9772	0.0254
8	31.70	0.9834	0.0156
9	32.70	0.9842	0.0162
10	34.01	0.9879	0.0149

570 E More Visualization Comparisons

571 We further compare DiffGS with DiffTF [4] on the airplane class of the ShapeNet dataset. The visual
 572 comparison is shown in Fig. 7, where our proposed DiffGS achieves more visually appealing results
 573 than DiffTF. DiffTF often produces shape generations with blurry textures, resulting in poor rendering
 574 quality. In contrast, our method produces high-fidelity renderings with the generated high-quality
 575 3DGS, accurately capturing both geometry and appearances.

576 F Limitation

577 One limitation of our method is that it sometimes produces overly creative color schemes. We
 578 illustrate this issue with two examples in Fig. 8. As shown, the failure cases of DiffGS result in
 579 excessively colorful appearances. For instance, the chair shown exhibits a color transition from



Figure 8: Failure cases of our method. DiffGS sometimes generates overly creative shapes with colorful appearances.

580 yellow to blue to green and finally to orange. These overly creative generations may not accurately
581 reflect real-world shapes.

582 **G Video**

583 We provide a video as part of our supplementary materials, which includes visualizations of the
584 unconditional generations, image/text conditional generations, Gaussian Splatting completion results,
585 and Point-to-Gaussian generations.

586 **NeurIPS Paper Checklist**

587 **1. Claims**

588 Question: Do the main claims made in the abstract and introduction accurately reflect the
589 paper's contributions and scope?

590 Answer: [Yes]

591 Justification: Our main claims made in the abstract and introduction accurately reflect the
592 paper's contributions and scope.

593 Guidelines:

- 594 • The answer NA means that the abstract and introduction do not include the claims
595 made in the paper.
- 596 • The abstract and/or introduction should clearly state the claims made, including the
597 contributions made in the paper and important assumptions and limitations. A No or
598 NA answer to this question will not be perceived well by the reviewers.
- 599 • The claims made should match theoretical and experimental results, and reflect how
600 much the results can be expected to generalize to other settings.
- 601 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
602 are not attained by the paper.

603 **2. Limitations**

604 Question: Does the paper discuss the limitations of the work performed by the authors?

605 Answer: [Yes]

606 Justification: We analysis the limitations of our method in Sec.F of the appendix.

607 Guidelines:

- 608 • The answer NA means that the paper has no limitation while the answer No means that
609 the paper has limitations, but those are not discussed in the paper.
- 610 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 611 • The paper should point out any strong assumptions and how robust the results are to
612 violations of these assumptions (e.g., independence assumptions, noiseless settings,
613 model well-specification, asymptotic approximations only holding locally). The authors
614 should reflect on how these assumptions might be violated in practice and what the
615 implications would be.
- 616 • The authors should reflect on the scope of the claims made, e.g., if the approach was
617 only tested on a few datasets or with a few runs. In general, empirical results often
618 depend on implicit assumptions, which should be articulated.
- 619 • The authors should reflect on the factors that influence the performance of the approach.
620 For example, a facial recognition algorithm may perform poorly when image resolution
621 is low or images are taken in low lighting. Or a speech-to-text system might not be
622 used reliably to provide closed captions for online lectures because it fails to handle
623 technical jargon.
- 624 • The authors should discuss the computational efficiency of the proposed algorithms
625 and how they scale with dataset size.
- 626 • If applicable, the authors should discuss possible limitations of their approach to
627 address problems of privacy and fairness.
- 628 • While the authors might fear that complete honesty about limitations might be used by
629 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
630 limitations that aren't acknowledged in the paper. The authors should use their best
631 judgment and recognize that individual actions in favor of transparency play an impor-
632 tant role in developing norms that preserve the integrity of the community. Reviewers
633 will be specifically instructed to not penalize honesty concerning limitations.

634 **3. Theory Assumptions and Proofs**

635 Question: For each theoretical result, does the paper provide the full set of assumptions and
636 a complete (and correct) proof?

637 Answer: [NA]

638 Justification: The paper does not include theoretical results.

639 Guidelines:

- 640 • The answer NA means that the paper does not include theoretical results.
- 641 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 643 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 644 • The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- 645 • Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- 646 • Theorems and Lemmas that the proof relies upon should be properly referenced.

650 4. Experimental Result Reproducibility

651 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
652 perimental results of the paper to the extent that it affects the main claims and/or conclusions
653 of the paper (regardless of whether the code and data are provided or not)?

654 Answer: [Yes]

655 Justification: We provide the detailed information in reproducing our methods in Sec.3,
656 Sec.4 of the main paper and the appendix. We also provide a demonstration code of our
657 method in the supplementary materials.

658 Guidelines:

- 659 • The answer NA means that the paper does not include experiments.
- 660 • If the paper includes experiments, a No answer to this question will not be perceived
661 well by the reviewers: Making the paper reproducible is important, regardless of
662 whether the code and data are provided or not.
- 663 • If the contribution is a dataset and/or model, the authors should describe the steps taken
664 to make their results reproducible or verifiable.
- 665 • Depending on the contribution, reproducibility can be accomplished in various ways.
666 For example, if the contribution is a novel architecture, describing the architecture fully
667 might suffice, or if the contribution is a specific model and empirical evaluation, it may
668 be necessary to either make it possible for others to replicate the model with the same
669 dataset, or provide access to the model. In general, releasing code and data is often
670 one good way to accomplish this, but reproducibility can also be provided via detailed
671 instructions for how to replicate the results, access to a hosted model (e.g., in the case
672 of a large language model), releasing of a model checkpoint, or other means that are
673 appropriate to the research performed.
- 674 • While NeurIPS does not require releasing code, the conference does require all submissions
675 to provide some reasonable avenue for reproducibility, which may depend on the
676 nature of the contribution. For example
 - 677 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
678 to reproduce that algorithm.
 - 679 (b) If the contribution is primarily a new model architecture, the paper should describe
680 the architecture clearly and fully.
 - 681 (c) If the contribution is a new model (e.g., a large language model), then there should
682 either be a way to access this model for reproducing the results or a way to reproduce
683 the model (e.g., with an open-source dataset or instructions for how to construct
684 the dataset).
 - 685 (d) We recognize that reproducibility may be tricky in some cases, in which case
686 authors are welcome to describe the particular way they provide for reproducibility.
687 In the case of closed-source models, it may be that access to the model is limited in
688 some way (e.g., to registered users), but it should be possible for other researchers
689 to have some path to reproducing or verifying the results.

690 5. Open access to data and code

691 Question: Does the paper provide open access to the data and code, with sufficient instruc-
692 tions to faithfully reproduce the main experimental results, as described in supplemental
693 material?

694 Answer: [Yes]

695 Justification: We provide our demonstration code as a part of our supplementary materials.
696 We will release the source code, data and instructions upon acceptance.

697 Guidelines:

- 698 • The answer NA means that paper does not include experiments requiring code.
- 699 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
700 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 701 • While we encourage the release of code and data, we understand that this might not be
702 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
703 including code, unless this is central to the contribution (e.g., for a new open-source
704 benchmark).
- 705 • The instructions should contain the exact command and environment needed to run to
706 reproduce the results. See the NeurIPS code and data submission guidelines ([https://nips.cc/
707 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 708 • The authors should provide instructions on data access and preparation, including how
709 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 710 • The authors should provide scripts to reproduce all experimental results for the new
711 proposed method and baselines. If only a subset of experiments are reproducible, they
712 should state which ones are omitted from the script and why.
- 713 • At submission time, to preserve anonymity, the authors should release anonymized
714 versions (if applicable).
- 715 • Providing as much information as possible in supplemental material (appended to the
716 paper) is recommended, but including URLs to data and code is permitted.

717 6. Experimental Setting/Details

718 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
719 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
720 results?

721 Answer: [Yes]

722 Justification: We provide the training and testing details in the experiment section (Sec.4)
723 and the appendix.

724 Guidelines:

- 725 • The answer NA means that the paper does not include experiments.
- 726 • The experimental setting should be presented in the core of the paper to a level of detail
727 that is necessary to appreciate the results and make sense of them.
- 728 • The full details can be provided either with the code, in appendix, or as supplemental
729 material.

730 7. Experiment Statistical Significance

731 Question: Does the paper report error bars suitably and correctly defined or other appropriate
732 information about the statistical significance of the experiments?

733 Answer: [No]

734 Justification: We report the average performance as the experimental results.

735 Guidelines:

- 736 • The answer NA means that the paper does not include experiments.
- 737 • The authors should answer “Yes” if the results are accompanied by error bars, confi-
738 dence intervals, or statistical significance tests, at least for the experiments that support
739 the main claims of the paper.
- 740 • The factors of variability that the error bars are capturing should be clearly stated (for
741 example, train/test split, initialization, random drawing of some parameter, or overall
742 run with given experimental conditions).

- 743 • The method for calculating the error bars should be explained (closed form formula,
 744 call to a library function, bootstrap, etc.)
 745 • The assumptions made should be given (e.g., Normally distributed errors).
 746 • It should be clear whether the error bar is the standard deviation or the standard error
 747 of the mean.
 748 • It is OK to report 1-sigma error bars, but one should state it. The authors should
 749 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
 750 of Normality of errors is not verified.
 751 • For asymmetric distributions, the authors should be careful not to show in tables or
 752 figures symmetric error bars that would yield results that are out of range (e.g. negative
 753 error rates).
 754 • If error bars are reported in tables or plots, The authors should explain in the text how
 755 they were calculated and reference the corresponding figures or tables in the text.

756 8. Experiments Compute Resources

757 Question: For each experiment, does the paper provide sufficient information on the com-
 758 puter resources (type of compute workers, memory, time of execution) needed to reproduce
 759 the experiments?

760 Answer: [Yes]

761 Justification: The computer resources needed to reproduce the experiments are provided in
 762 Sec.C of the appendix.

763 Guidelines:

- 764 • The answer NA means that the paper does not include experiments.
- 765 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
 766 or cloud provider, including relevant memory and storage.
- 767 • The paper should provide the amount of compute required for each of the individual
 768 experimental runs as well as estimate the total compute.
- 769 • The paper should disclose whether the full research project required more compute
 770 than the experiments reported in the paper (e.g., preliminary or failed experiments that
 771 didn't make it into the paper).

772 9. Code Of Ethics

773 Question: Does the research conducted in the paper conform, in every respect, with the
 774 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

775 Answer: [Yes]

776 Justification: The research conducted in the paper conform, in every respect, with the
 777 NeurIPS Code of Ethics.

778 Guidelines:

- 779 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 780 • If the authors answer No, they should explain the special circumstances that require a
 781 deviation from the Code of Ethics.
- 782 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
 783 eration due to laws or regulations in their jurisdiction).

784 10. Broader Impacts

785 Question: Does the paper discuss both potential positive societal impacts and negative
 786 societal impacts of the work performed?

787 Answer: [Yes]

788 Justification: We discuss the applications and potential impacts of our method in the
 789 introduction.

790 Guidelines:

- 791 • The answer NA means that there is no societal impact of the work performed.
- 792 • If the authors answer NA or No, they should explain why their work has no societal
 793 impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We use the open-sourced datasets under their licenses.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- 848 • For existing datasets that are re-packaged, both the original license and the license of
849 the derived asset (if it has changed) should be provided.
850 • If this information is not available online, the authors are encouraged to reach out to
851 the asset's creators.

852 **13. New Assets**

853 Question: Are new assets introduced in the paper well documented and is the documentation
854 provided alongside the assets?

855 Answer: [NA]

856 Justification: The paper does not release new assets.

857 Guidelines:

- 858 • The answer NA means that the paper does not release new assets.
859 • Researchers should communicate the details of the dataset/code/model as part of their
860 submissions via structured templates. This includes details about training, license,
861 limitations, etc.
862 • The paper should discuss whether and how consent was obtained from people whose
863 asset is used.
864 • At submission time, remember to anonymize your assets (if applicable). You can either
865 create an anonymized URL or include an anonymized zip file.

866 **14. Crowdsourcing and Research with Human Subjects**

867 Question: For crowdsourcing experiments and research with human subjects, does the paper
868 include the full text of instructions given to participants and screenshots, if applicable, as
869 well as details about compensation (if any)?

870 Answer: [NA]

871 Justification: The paper does not involve crowdsourcing nor research with human subjects.

872 Guidelines:

- 873 • The answer NA means that the paper does not involve crowdsourcing nor research with
874 human subjects.
875 • Including this information in the supplemental material is fine, but if the main contribu-
876 tion of the paper involves human subjects, then as much detail as possible should be
877 included in the main paper.
878 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
879 or other labor should be paid at least the minimum wage in the country of the data
880 collector.

881 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
882 Subjects**

883 Question: Does the paper describe potential risks incurred by study participants, whether
884 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
885 approvals (or an equivalent approval/review based on the requirements of your country or
886 institution) were obtained?

887 Answer: [NA]

888 Justification: The paper does not involve crowdsourcing nor research with human subjects.

889 Guidelines:

- 890 • The answer NA means that the paper does not involve crowdsourcing nor research with
891 human subjects.
892 • Depending on the country in which research is conducted, IRB approval (or equivalent)
893 may be required for any human subjects research. If you obtained IRB approval, you
894 should clearly state this in the paper.
895 • We recognize that the procedures for this may vary significantly between institutions
896 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
897 guidelines for their institution.
898 • For initial submissions, do not include any information that would break anonymity (if
899 applicable), such as the institution conducting the review.