

GAP: Gaussianize Any Point Clouds with Text Guidance

Supplementary Material

001 1. Source Codes

002 We provide our demonstration code as a part of our supple-
003 mentary materials. We will release the source code, data,
004 and instructions upon acceptance.

005 2. More Implementation Details

006 Our text-to-image generation pipeline is based on the Sta-
007 ble Diffusion v1.5 model [4] with ControlNet [7] for depth
008 conditioning. Specifically, we utilize the Depth-to-Image
009 model architecture that incorporates depth maps as geomet-
010 ric guidance during the diffusion process. For Gaussian op-
011 timization, each viewpoint undergoes 1000 iterations of op-
012 timization. During optimization, we set the distance loss
013 weight α to 1.5 and the scale loss weight β to 5. For the
014 Scale Loss, we set the maximum scale threshold τ to $1e-6$.
015 In the spatial-aware Gaussian inpainting stage, we set the
016 number of nearest neighbors $L = 90$ for color diffusion and
017 the radius $\rho = 0.1$ for opacity control. The base opacity
018 value o_0 is set to 5 and the density threshold P_0 is set to
019 100. All experiments are conducted on NVIDIA RTX3090.
020 The complete pipeline for processing an object takes about
021 25 minutes.

022 **Adaption of Baselines in Point-to-Gaussian Genera-**
023 **tion.** We extended DreamGaussian, originally designed for
024 image/text-to-3D generation, to support point cloud inputs
025 by replacing its random Gaussian initialization with point
026 cloud-guided initialization. We adapt TriplaneGaussian as
027 our baseline by modifying its original image-conditioned
028 3DGS generation pipeline. Specifically, we bypass its point
029 cloud decoder for direct point-to-Gaussian conversion and
030 incorporate Stable Diffusion to enable text-to-image gen-
031 eration. DiffGS uses a Gaussian VAE to convert point clouds
032 into Gaussians by querying features from triplanes. How-
033 ever, its lack of text-guided appearance control limited us to
034 conducting only unconditional point-to-gaussian generation
035 experiments.

036 3. Evaluation Metrics

037 We employ three complementary metrics to comprehen-
038 sively evaluate our method’s performance: Fréchet Incep-
039 tion Distance (FID) [6], Kernel Inception Distance (KID)
040 [2], and CLIP Score [3]. For each metric, we render the
041 generated results from fixed viewpoints at 1024×1024 res-
042 olution.

043 3.1. Fréchet Inception Distance (FID)

044 FID measures the similarity between the distribution of gen-
045 erated images and real images. We compute feature repre-
046 sentations using the InceptionV3 [5] network pretrained on
047 ImageNet. The FID score is calculated as:

$$\text{FID} = \|\mu_r - \mu_g\|^2 + \text{Tr} \left(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{1/2} \right), \quad (1) \quad 048$$

049 where μ_r , Σ_r and μ_g , Σ_g are the mean and covariance
050 matrices of the real and generated feature distributions re-
051 spectively. Lower FID scores indicate better generation
052 quality.

053 3.2. Kernel Inception Distance (KID)

054 KID provides an unbiased estimate of the Maximum Mean
055 Discrepancy (MMD) between real and generated image fea-
056 tures. We report KID scores multiplied by 10^3 for better
057 readability. The KID metric is computed as:

$$\begin{aligned} \text{KID} &= \text{MMD}^2(X_r, X_g) \\ &= \frac{1}{n(n-1)} \sum_{i \neq j} k(x_i, x_j) + \frac{1}{m(m-1)} \sum_{i \neq j} k(y_i, y_j) \\ &\quad - \frac{2}{mn} \sum_{i,j} k(x_i, y_j), \end{aligned} \quad (2) \quad 058$$

059 where X_r and X_g are real and generated feature sets re-
060 spectively, and $k(\cdot, \cdot)$ is the polynomial kernel. Like FID,
061 lower KID scores indicate better quality.

062 3.3. CLIP Score

063 CLIP Score evaluates the semantic alignment between the
064 images rendered from Gaussians and the input text prompts.
065 We use the CLIP ViT-L/14 [3] model to compute the cosine
066 similarity between text and image embeddings. For each
067 generated result, we average the CLIP scores across all ren-
068 dered views. Higher CLIP scores indicate better text-image
069 alignment.

070 4. More Results

071 4.1. Text-Driven Appearance Generation

072 To demonstrate GAP’s advantage over mesh-based meth-
073 ods with reconstructed geometries, we provide extensive
074 visual comparisons in Fig. 2. We compare our results
075 with both traditional geometry-based reconstruction using

076 Ball-Pivoting Algorithm (BPA) [1] and learning-based re-
 077 construction using CAPUDF [8]. For each reconstruc-
 078 tion method, we generate UV maps using xatlas and apply
 079 the same texture generation methods (Texture, Text2Tex,
 080 Paint3D, SyncMVD) as baselines.

081 The visual comparison reveals two major challenges
 082 when using reconstructed meshes. First, both BPA [1] and
 083 CAPUDF [8] reconstructed meshes suffer from geom-
 084 etric ambiguities and information loss during surface recon-
 085 struction. Second, the excessive number of faces in recon-
 086 structed meshes leads to highly fragmented UV layouts with
 087 severe stretching and overlapping issues, as shown in Fig. 1.
 088 These fragmented UV charts not only limit the effective tex-
 089 ture resolution but also cause color bleeding artifacts across
 090 chart boundaries, resulting in discontinuities and inconsis-
 091 tencies in the final appearance.

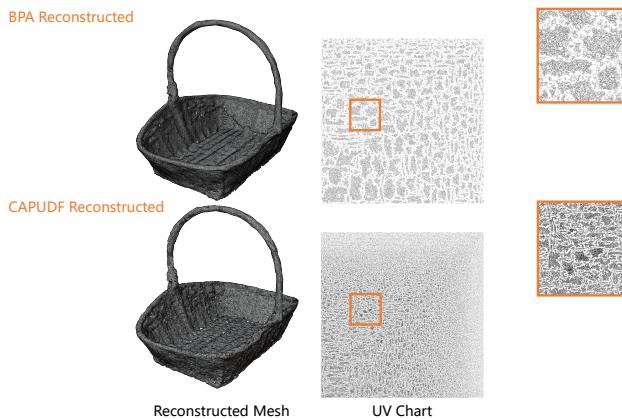


Figure 1. Comparison of UV parameterization results for reconstructed meshes using BPA and CAPUDF. The UV layouts exhibit severe fragmentation, stretching, and overlapping issues.

092 In contrast, GAP directly optimizes Gaussian primitives
 093 in 3D space without requiring intermediate mesh recon-
 094 struction or UV mapping. This direct optimization ap-
 095 proach preserves geometric details from the input point
 096 cloud while enabling high-quality appearance generation
 097 across different object categories.

098 4.2. More Application Results

099 We further show more visualizations and comparisons of
 100 the application shown in the main Paper. We show more
 101 comparisons on the task of Point-to-Gaussian generation in
 102 Fig. 3. More visualizations on the Gaussian generations un-
 103 der real-world scanned DeepFashion3D dataset is shown in
 104 Fig. 4. Finally, we show more comparisons on learning to
 105 generate appearances for real-world 3D scenes in Fig. 5.

5. Video

We provide a supplementary video demonstrating our results across different tasks, including text-guided appearance generation, point-to-Gaussian conversion, Gaussian synthesis from real-world scans, and scene-level Gaussian generation.

References

- [1] Fausto Bernardini, Joshua Mittleman, Holly Rushmeier, Cláudio Silva, and Gabriel Taubin. The ball-pivoting algorithm for surface reconstruction. *IEEE Transactions on Visualization and Computer Graphics*, 5(4):349–359, 1999. 2
- [2] Mikołaj Bińkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. *arXiv preprint arXiv:1801.01401*, 2018. 1
- [3] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 1
- [4] 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. 1
- [5] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2818–2826, 2016. 1
- [6] Yu Yu, Weibin Zhang, and Yun Deng. Frechet inception distance (fid) for evaluating gans. *China University of Mining Technology Beijing Graduate School*, 3, 2021. 1
- [7] X. Zhang, Y. Liu, and L. Zhang. Controlnet: Learning to control diffusion models with pretrained optimizers. *arXiv:2302.05543*, 2023. 1
- [8] Junsheng Zhou, Baorui Ma, Yu-Shen Liu, Yi Fang, and Zhizhong Han. Learning consistency-aware unsigned distance functions progressively from raw point clouds. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2

BPA Reconstruct*CAPUDF Reconstruct*

Figure 2. Visual comparison of text-guided appearance generation results with the reconstructed meshes.



Figure 3. Visual comparison of point-to-Gaussian generation results on ShapeNet Chair and DeepFashion3D. Our GAP method demonstrates superior visual quality and geometric accuracy with flexible text-guided appearance control.



Figure 4. Gaussianization results on real-world partial scans from SRB and DeepFashion3D datasets. Our surface-anchoring mechanism and diffusion-guided rendering supervision enable GAP to generate complete, high-quality 3D Gaussian representations while maintaining geometric consistency.

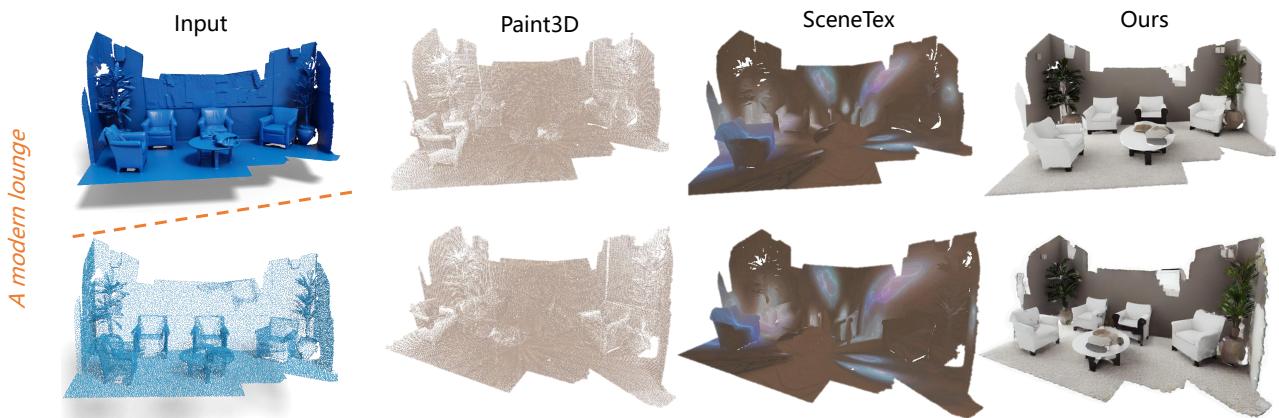


Figure 5. Scene-level Gaussianization comparison on real-world scanned 3DScene datasets. GAP generates high-quality results for complex scenes through a single optimization process.