

FocalDreamer: Text-Driven 3D Editing via Focal-Fusion Assembly

Yuhan Li¹, Yishun Dou², Yue Shi¹, Yu Lei¹, Xuanhong Chen¹, Yi Zhang¹, Peng Zhou¹, Bingbing Ni^{1†}

¹Shanghai Jiao Tong University, Shanghai 200240, China

²Huawei

{melodious, nibingbing}@sjtu.edu.cn

Abstract

While text-3D editing has made significant strides in leveraging score distillation sampling, emerging approaches still fall short in delivering separable, precise and consistent outcomes that are vital to content creation. In response, we introduce FocalDreamer, a framework that merges base shape with editable parts according to text prompts for fine-grained editing within desired regions. Specifically, equipped with geometry union and dual-path rendering, FocalDreamer assembles independent 3D parts into a complete object, tailored for convenient instance reuse and part-wise control. We propose geometric focal loss and style consistency regularization, which encourage focal fusion and congruent overall appearance. Furthermore, FocalDreamer generates high-fidelity geometry and PBR textures which are compatible with widely-used graphics engines. Extensive experiments have highlighted the superior editing capabilities of FocalDreamer in both quantitative and qualitative evaluations.

1 Introduction

Art reflects the figments of human imagination and creativity. Recently, the rapid development of neural generative models (Dhariwal and Nichol 2021) has significantly lowered the barriers for humans to engage in artistic creation with just a few words. However, these black-box models also deprive humans of a significant portion of control, which means the generation isn't often aligned with expectations. In this work, we take a step towards precise editing for 3D creation, enabling networks to naturally expand user's intentions, rather than controlling the entire generative process.

In the realms of animation, gaming, and the recent advance of virtual augmented reality, 3D models and scenes are commonly constructed as an assembly of semantically distinct base parts, which support the practice of rendering multiple copies of the same part across scenes with different transform matrices, called *geometry instancing* or *instance reuse* (Fig. 1). We believe that an ideal 3D editing workflow should possess the following good properties:

- **Separable.** Given a base shape, it should produce structurally *separate parts* (Li, Niu, and Xu 2020) facili-

Copyright © 2024, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

[†]Corresponding author: Bingbing Ni.

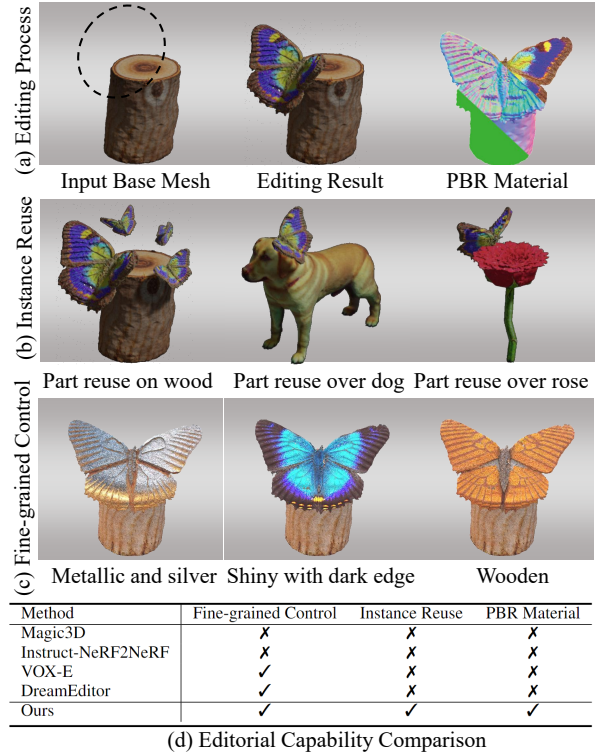


Figure 1: Given the prompt “a butterfly over a tree stump”, our method delivers high-fidelity geometry and photorealistic appearance using PBR materials. Lines (b-c) showcase FocalDreamer’s capability for separable and precise edits.

tating for instance reuse and part-wise post-processing, grounded in widespread understanding.

- **Precise.** It should provide *fine-grained* and *local* editing, enabling precise control in the desired area (Zhuang et al. 2023), while maintaining other regions untouched.
- **Consistent.** After the editing process, the resultant shape should respect the characteristics of the source shape in *harmonious appearance* (Xie et al. 2023), while visually adhering to the text specifications.

Emerging approaches in text-3D editing have achieved noteworthy development, yet they often fall short in deliv-

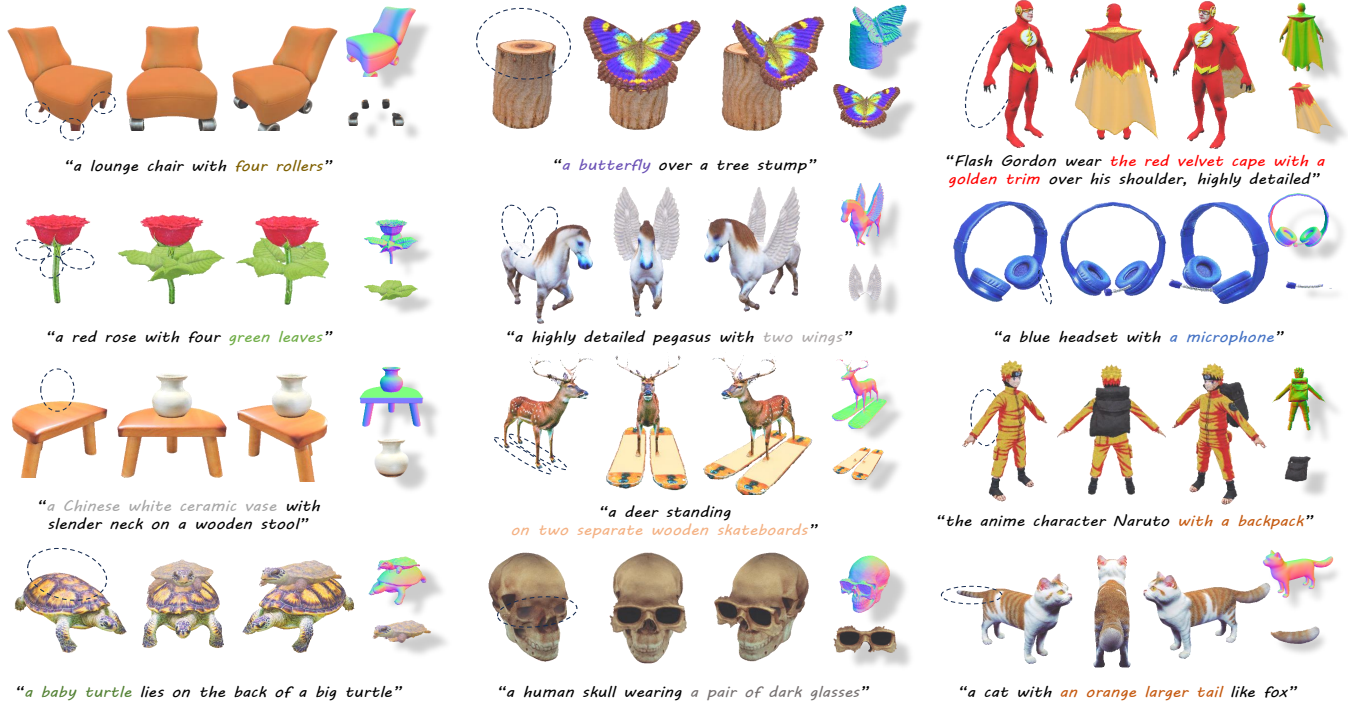


Figure 2: FocalDreamer can generate meticulously detailed and photo-realistic 3D editing. The left column displays base meshes with focal regions. The three right columns showcase edited overall appearance, assembled geometry, and editable part.

ering separable, precise, and consistent outcomes that are vital to content creation. Some approaches (Lin et al. 2023; Haque et al. 2023) struggle to pinpoint the focused local regions, leading to undesired alterations to the base shape. Others (Sella et al. 2023; Zhuang et al. 2023) overlook the stylistic consistency of the 3D edited portions. Furthermore, nearly all past methods directly modify base shape, neglecting the need for *instance reuse* and *part-wise control* (i.e., enabling fine-grained edits to individual parts of an object).

We introduce the following key contributions to meet our outlined criteria: (1) **Separable**: we propose FocalDreamer, a user-friendly framework that permits intuitive object modifications using text prompts and a rough focal region for the intended edits. Instead of direct modifications to the **base shape** (e.g., the *horse* in Fig. 3), a novel **editable part** (*wings* in Fig. 3) is generated in the focal region, facilitating instance reuse and precise control. Equipped with geometry union and dual-path rendering, this part is merged with base mesh into a semantically unified shape in a lossless and differentiable manner, then optimized using a powerful text-to-image model to align the prompts and shapes. Furthermore, our decoupled learning of geometry and appearance yields detailed geometry and PBR textures, ensuring compatibility with prominent graphics engines. (2) **Precise**: Users delineate one or several ellipsoid focal regions, in which a spherical editable part initializes, acting as a smooth prior for the geometry network. The geometric focal loss is also introduced, discouraging edits beyond specified regions. (3) **Consistent**: a smooth, coherent surface is essential in certain scenes. Hence, a soft geometry union operator and a style

consistency regularization are proposed to ensure a seamless geometric transition and stylistically consistent texture between the learnable part and base shape.

To our knowledge, this is the first component-based editing method with separate learnable parts. Rich experiments and detailed ablation studies highlight the superior editing capabilities of our approach, as shown in Fig. 2.

2 Related Work

Text-guided Image Generation and Editing. Significant progress in Text-to-Image (T2I) generation with diffusion models (Ho, Jain, and Abbeel 2020) is witnessed in recent years. More recently, with the availability of scalable generator architectures and extremely large-scale image-text paired datasets, they’ve demonstrated impressive performance in high-fidelity and flexible image synthesis (Romach et al. 2022). Due to their comprehension of complex concepts, diffusion models are also amicable for various editing tasks, such as image inpainting (Lugmayr et al. 2022), image stylization (Zhang et al. 2023). The most relevant field to us among those is inpainting, which provides flexible control of the inpainted content, and a mask to constrain the shape of the inpainted object. SmartBrush (Xie et al. 2023) introduces a precision factor into the masks for multiple-grained controls on inpainting regions.

Text-to-3D Content Generation. Driven by the aspiration to produce high-fidelity 3D content using semantic inputs like text prompts, the field of text-to-3D has garnered a significant boost in recent years (Poole et al. 2022). Ear-

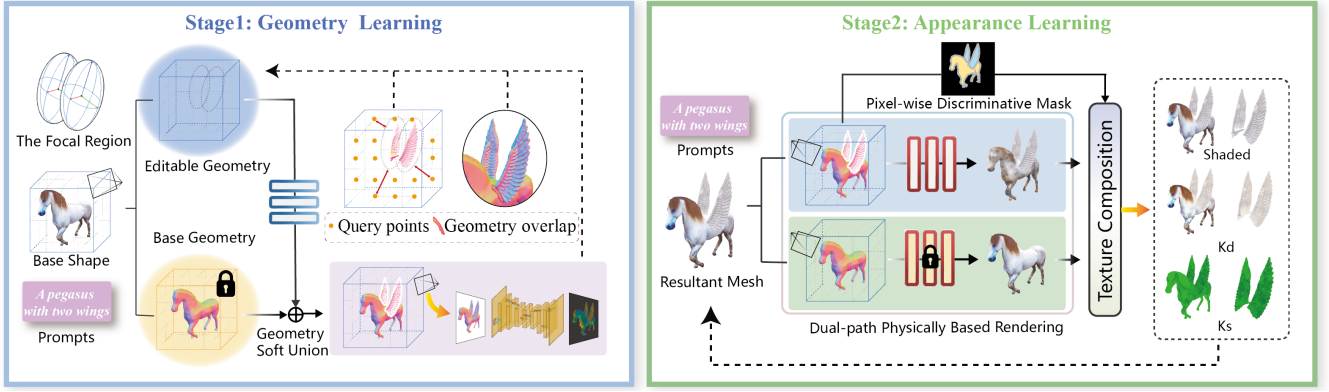


Figure 3: An overview of FocalDreamer. (a) During geometry learning, given a base shape, we first initialize an ellipsoid as editable geometry within each focal region. Then we render the normal map of merged shape as shape encoding of pre-trained T2I models, to optimize the editable geometry according to prompts. (b) During appearance learning, resultant shape is rendered in a dual-path manner with base and editable textures. The outcomes are then blended by Pixel-wise Discriminative Mask for a unified appearance. (c) Several regularizations are introduced to improve the editing quality, including \mathcal{L}_{GF} , \mathcal{L}_{CA} , and \mathcal{L}_{SC} .

lier approaches either align shapes and images in the latent space by CLIP supervision (Radford et al. 2021) to generate 3D geometries (Mohammad Khalid et al. 2022) or synthesize new perspectives (Jain et al. 2022), or they train text-conditioned 3D generative models from the ground up (Li et al. 2023). DreamFusion (Poole et al. 2022) first employs large-scale T2I models with a combination of score distillation sampling to distill the prior, and achieves impressive results. Magic3D (Lin et al. 2023) further improved the quality and performance of generated 3D shapes with a 2-step pipeline. TextMesh (Tsalicoglou et al. 2023) modify the 3D representation to extract detailed mesh. However, all these methods present semantic misalignment between the local content and global text description when editing, leaning towards distorted background and inconsistent results.

3D Content Editing. Semantic-driven 3D scene editing is a much harder task compared with 2D photo editing because of the high demand for multi-view consistency, the scarcity of paired 3D data and its entangled geometry and appearance. Previous approaches either rely on laborious annotation (Kania et al. 2022; Yang et al. 2022), only support object deformation or translation (Tschernezki et al. 2022; Kobayashi, Matsumoto, and Sitzmann 2022), or only perform global style transfer (Chen et al. 2022; Chiang et al. 2022; Fan et al. 2022; Huang et al. 2022) without strong semantic meaning. Recently, thanks to the development of score distillation sampling technique, text-guided editing has emerged as a promising direction with great potential. SKED (Mikaeili et al. 2023) possesses the capability to edit 3D scenes with multi-view sketches. Latent-NeRF (Metzer et al. 2023) and Fantasia3D (Chen et al. 2023) realize sketch-shape guidance by relaxed geometric constraints. Instruct-NeRF2NeRF (Haque et al. 2023) can edit an existing NeRF scene by iterative dataset update. However, it manipulates the entire space, and the preservation of undesired regions is absent. Vox-E (Sella et al. 2023) allows local edits on an existing NeRF, but it suffers from subpar editing quality and

noticeable noise as shown in Section 4, because of coupling geometry and textures. Most related to our work, DreamEditor (Zhuang et al. 2023) locally edits a mesh-based neural field. However, it doesn’t achieve separable editing which is vital for instance reuse and part-wise control. Moreover, DreamEditor cannot change the number of vertices, supporting only minor shape insertion and replacement of objects of the same type (e.g., a horse to a deer). In contrast, our work not only brings about reasonable and noticeable geometric changes but also generates realistic appearances.

3 Method

As illustrated in Fig. 3, a complete object is conceptualized as a composition of base shape and learnable parts, wherein both of them possess their own geometry and texture, tailored for convenient instance reuse and part-wise control. Furthermore, a two-stage training strategy is adopted to sequentially learn the geometry and texture of the editable shape, to avoid the potential interference that can occur when geometry and texture learning are intertwined. For instance, in the case of *zebra* modeling, geometric protrusions might be learned instead of the desired black stripes. Such a disentangled representation not only stabilizes the training process but also yields high-fidelity geometry and textures, especially when compared to popular text-to-3D models.

3.1 Preliminary

Score Distillation Sampling. Score distillation sampling (SDS) is a way to distill the priors hidden in large T2I models for 3D generation proposed by DreamFusion (Poole et al. 2022). DreamFusion represents 3D scenes as a series of learnable parameters θ . Utilizing a differentiable renderer, it converts the 3D scenes into 2D image sets x . Subsequently, it employs large-scale models ϕ to optimize the parameters of the 3D scenes with a score function as follows:

$$\nabla_{\theta} \mathcal{L}_{SDS}(\phi, x) = \mathbb{E}_{t, \epsilon} \left[w(t) (\hat{e}_{\phi}(z_t; y, t) - \epsilon) \frac{\partial x}{\partial \theta} \right], \quad (1)$$

where $w(t)$ controls the weight of SDS guidance depending on noise level t . $\hat{\epsilon}_\phi(z_t; y, t)$ and ϵ are the predicted noise and actual noise, respectively. y is the condition.

DMTet. DMTet (Munkberg et al. 2022) is a hybrid representation that has two components, *i.e.*, a deformable tetrahedral grid and a differentiable Marching Tetrahedral (MT) layer. The Signed Distance Function (SDF) values and the position offsets of deformable tetrahedral vertices are learnable, followed by the MT layer to extract meshes.

3.2 Geometry Editing

Focal Region. The starting point of our algorithm is a base shape (Ψ_b for geometry and Γ_b for texture) to be edited, which can be the reconstruction from images, crafted shapes by artists (Munkberg et al. 2022), and even the novel shapes from the generative method (Chen et al. 2023). Then the base model is modified by compositing with a new learnable part according to prompts. To offer more precise control over the generation process, users are requested to select one or multiple ellipsoid areas (depending on the editing needs) as focal/target regions. Each focal region Ω' is deformed from a standard sphere Ω by an affine transformation with 9 degrees of freedom (DOF), 3 DOF for stretching, 3 DOF for rotation, and 3 DOF for translation along the $\{X, Y, Z\}$ -axis:

$$\Omega' = R_{xyz}(\alpha, \beta, \gamma) \cdot T(t_x, t_y, t_z) \cdot S(s_x, s_y, s_z) \cdot \Omega. \quad (2)$$

The selection of the focal region doesn't require exact precision for it merely serves as a rough expression of the regional prior from user intent. Our model will optimally generate geometry driven by the text input. Furthermore, we initialize ellipsoids within specified regions, offering a smooth prior that enhances the stability of the geometric modeling.

Geometry Learning and Fusion. We adopt DMTet as our 3D scene representation optimized by the prior knowledge distilled from pre-trained T2I model. More specifically, keeping the base shape $\Psi_b(v_i)$ frozen, we parameterize the SDF values (inner is positive) of editable parts using MLP $\Psi_e(v_i)$ for each vertex v_i within the tetrahedral grid. Subsequently, a soft geometry union (Quilez and Jeremias 2018) is performed between $\Psi_b(v_i)$ and $\Psi_e(v_i)$, resulting in $\Psi_u(v_i)$ for a smooth junction:

$$\Psi_u(v_i) = \max\{\Psi_b(v_i), \Psi_e(v_i)\} + \frac{0.1 \times h^2}{k}, \quad (3)$$

$$\text{where } h = \max\{(k - |\Psi_b(v_i) - \Psi_e(v_i)|), 0\}, \quad (4)$$

where k determines the extent of the soft merge and is set to 0.15 by default. After geometry fusion, a differentiable MT layer transforms $\Psi_u(v_i)$ and the vertex offset Δv_i into a triangular surface mesh \mathcal{M} . Finally, the rendered normal map n and the object mask o extracted from the mesh \mathcal{M} are fed into pre-trained T2I models with SDS loss to update Ψ_e :

$$\nabla_{\Psi_e} \mathcal{L}_{SDS}(\phi, \tilde{n}) = \mathbb{E}_{t, \epsilon} \left[w(t)(\hat{\epsilon}_\phi(z_t^{\tilde{n}}; y, t) - \epsilon) \frac{\partial \tilde{n}}{\partial \Psi} \frac{\partial z^{\tilde{n}}}{\partial \tilde{n}} \right], \quad (5)$$

where ϕ parameterized pre-train T2I model, \tilde{n} represents the augmentation of n concatenated with o , $z^{\tilde{n}}$ is latent encoding of \tilde{n} . We observed using normal map n promotes the

expression of geometric details and training stability (Chen et al. 2023). This improvement from n is partly attributed to disentangling the geometry from the intertwining of texture, and its sufficient expressiveness to depict complex geometric details.

Geometric Concentration. One of the main criteria for a proficient 3D editing algorithm is its ability to retain the geometry and color of the base object throughout the editing process. However, the aforementioned pipeline cannot ensure locality in editing. We have observed global changes and a loss of characteristics from the base shape (Fig. 7). To counteract it, we introduce distance-aware **geometric focal loss** \mathcal{L}_{GF} . During each iteration, a certain number of points $p_i \in \mathbb{R}^3$ are sampled outside the user-specified focal region Ω' , with their SDF values $\Psi_e(p_i)$ and their distances d_i to the focal region Ω' . The objective of \mathcal{L}_{GF} is punishing the editable shape when it produces topological structures ($\Psi_e(p_i) > 0$) outside Ω' . Moreover, the closer p_i is to the target region, the less the penalty, for this distance-aware setting permits geometry to overrun beyond the rough focal region slightly. The **geometric focal loss** is defined as:

$$\mathcal{L}_{GF} = \mathbb{E}_{p_i \notin \Omega'} \left[\left(1 - e^{\frac{-d_i^2}{\sigma_1}} \right) \cdot \tanh\left(\frac{\max\{\Psi_e(p_i) + \xi, 0\}}{\sigma_2}\right) \right], \quad (6)$$

where $\sigma_1 = 0.05$ and $\sigma_2 = 0.01$ control how sensitive the loss is, *i.e.*, lower $\sigma_{1,2}$ values tighten the constraint on the optimization such that only the editable region is modified strictly. The hyperparameter ξ is a small positive threshold to prevent topological structures from minor positive SDF values. For computational efficiency, we sample query points on the tetrahedral vertex v_i , and pre-compute their distance d_i to Ω' before the geometry generation process begins.

Collision Avoidance. Another essential criterion is to respect the purity of the editing results, *i.e.*, the editable shape should not overlap with the base shape, as they are semantically independent and distinct parts. We enforce it by penalizing the query points p_i that reside both within the learnable shape and the base shape with the **collision avoidance loss**:

$$\mathcal{L}_{CA} = \mathbb{E}_{p_i} [\max\{\Psi_b(p_i), 0\} \cdot \max\{\Psi_e(p_i), 0\}]. \quad (7)$$

Intuitively, this reduces the likelihood of overlap between the editable shape and the original mesh, resulting in cleaner editing outcomes. For computational efficiency, we sample query points at v_i as the same as geometric focal Loss.

3.3 Appearance Editing

Dual-path Physically Based Rendering. After the optimization of the geometry network, the resultant mesh \mathcal{M} is obtained from the soft fusion and MT layer. Following Physically Based Rendering (PBR) material model, we use hash-grid-based texture neural fields Γ for \mathcal{M} to produce the diffuse term k_d , the roughness and metallic term k_{rm} , and the normal term k_n as $(k_d, k_{rm}, k_n) = \Gamma(p_i)$. In order to retain the appearance of the base shape untouched, a naive and straightforward idea would be to initialize the learnable texture neural fields Γ_e with the base texture fields

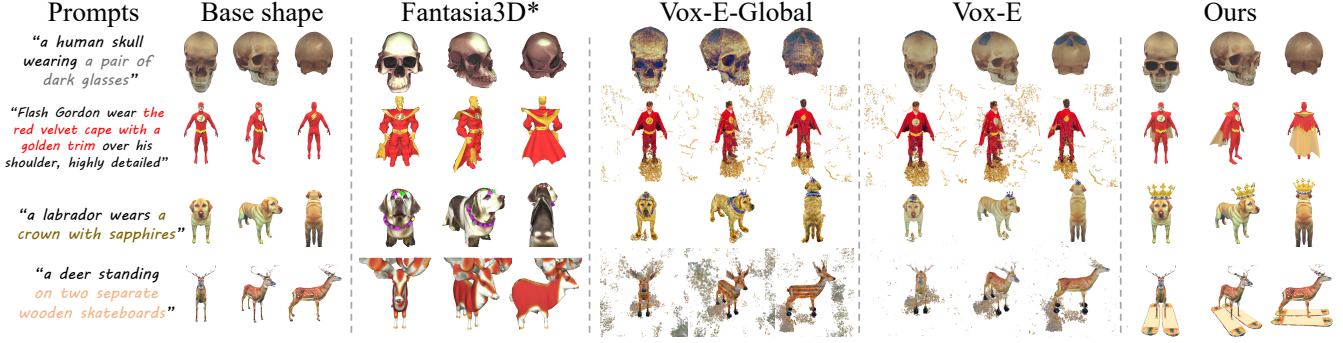


Figure 4: Visual comparison. Our approach synthesizes high-quality edits while preserving the base mesh perfectly.

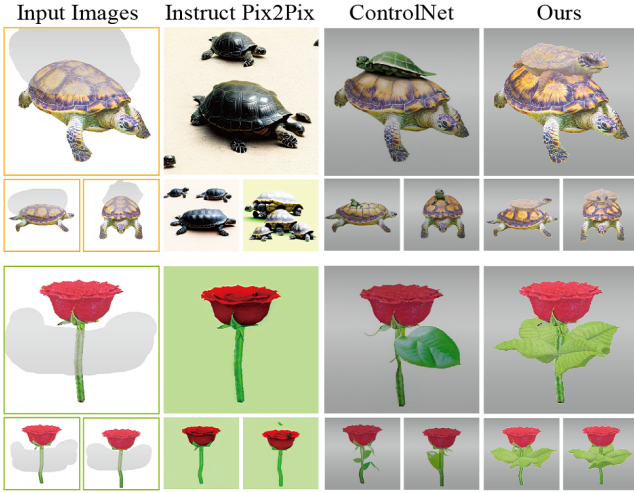


Figure 5: Comparison with SOTA image editing methods. The gray areas in input images indicate the in-painting region. We observed that 2D editing methods exhibit view-inconsistent, and their quality varies with viewpoints.

Method	CLIP _{sim} ↑	CLIP _{dir} ↑
Fantasia3D*	0.284	0.0180
Vox-E-Global	0.299	0.0204
Vox-E	0.293	0.0178
FocalDreamer (ours)	0.329	0.0519

Table 1: Quantitative evaluation results across 15 scenes.

Γ_b derived from the base shape reconstruction, then the entire shape’s appearance is modeled by Γ_e exclusively. However, this simple pipeline has two shortcomings: 1) As the number of iterations increases, it suffers from sub-optimal convergence and loss of the original material (in Fig. 7). In essence, the texture of the base shape isn’t adequately retained due to the overly strong knowledge supervision from T2I models. 2) Although learnable parts have independent semantics, such as “the wings”, their texture cannot be extracted alone. This impediment makes the reuse and driving of materials for these editable parts unfeasible.

To tackle this issue, we re-design the rendering pipeline

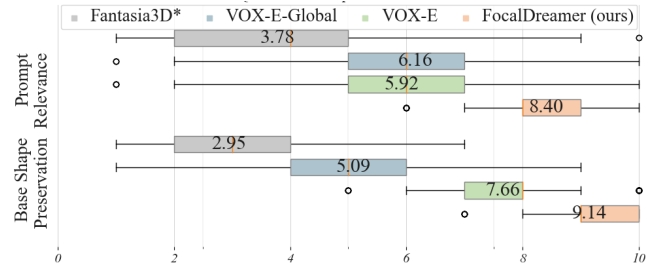


Figure 6: Boxplot illustration of user study. FocalDreamer demonstrates better performance (high means) and stability across scenes (narrow interquartile range).

in a dual-path manner. Central to this redesign is a Pixel-wise Discriminative Mask (PDM) generated in the rasterization process, which discerns whether each pixel comes from the face of the base mesh or the editable mesh. As depicted in Fig. 3, throughout the dual-path rendering process, both parts are rendered based on their own neural texture fields, and the outcomes are then blended by PDM, which is called texture composition, culminating in a unified merged view. Similarly, the merged view is inputted into the T2I model for texture optimization with SDS loss. By truncating the gradient towards Γ_b , the texture of the base shape is precisely preserved, while the editable shape has its independent trainable texture Γ_e . Dual-path rendering balances the preservation of the base shape structure with flexible part-wise control, as well as the seamless integration of both parts.

Style Consistency. In some instances, local changes are anticipated to be realized seamlessly, as well as in a harmoniously coordinated style, as shown in Fig. 7. This problem is modeled as follows: let $\mathcal{M}_e \in \mathbb{R}^3$ be a closed subspace to represent the editable parts with boundary $\partial\mathcal{M}_e$. Let f^* be a known mapping function defined over \mathbb{R}^3 minus the interior of \mathcal{M}_e to be preserved, and let f be the unknown function defined over the interior of \mathcal{M}_e . A classical interpolant f is defined as the solution (Pérez, Gangnet, and Blake 2003):

$$\min_f \iint_{\mathcal{M}_e} |\nabla f|^2 \text{ with } f|_{\partial\mathcal{M}_e} = f^*|_{\partial\mathcal{M}_e}. \quad (8)$$

We propose two consistency regularization items to imi-

\mathcal{L}_{GF}	\mathcal{L}_{CA}	\mathcal{L}_{SC}	Dual-path Render	CLIP _{sim}	CLIP _{dir}
✓	✓	✗	✓	0.312*	0.0402*
✓	✓	✓	✓	0.319*	0.0495*
✗	✓	✓	✓	0.316	0.0433
✓	✗	✓	✓	0.329	0.0517
✓	✓	✓	✗	0.313	0.0401
✓	✓	✓	✓	0.329	0.0519

Table 2: Ablation study. Since not all scenes require style consistency, we report the editings require \mathcal{L}_{SC} with *.

tate the interpolant process :

$$\mathcal{L}_g = \mathbb{E}_{p_i \in \mathcal{M}_e} \left[\|\Gamma_e(p_i) - \Gamma_e(p_i + \delta)\|^2 \right], \quad (9)$$

$$\mathcal{L}_b = \mathbb{E}_{p_i \in \partial \mathcal{M}_e} \left[\|\Gamma_e(p_i) - \Gamma_b(p_i)\|^2 \right], \quad (10)$$

$$\mathcal{L}_{SC} = \mathcal{L}_g + \lambda \mathcal{L}_b. \quad (11)$$

Intuitively, the \mathcal{L}_b ensures that the editable texture Γ_e is consistent with the base texture Γ_b in the adjoining areas $\partial \mathcal{M}_e$ as Dirichlet boundary condition, while the \mathcal{L}_g extends the consistent style throughout the whole learnable part Γ_e with gradient constrain on small noise δ .

4 Experiments

4.1 Experimental Setups

Implementation Details. We use the Stable Diffusion implementation by HuggingFace Diffusers for SDS, and adopt DMTet to learn geometry and texture separately with NVDiffRast as a differentiable renderer. FocalDreamer usually takes less than 30 minutes (3000 steps) for geometry and 20 minutes (2000 steps) for texture to converge on 4 Nvidia RTX 3090 GPUs, where we use AdamW optimizer with the respective learning rates of 1×10^{-3} and 1×10^{-2} . UV edge padding techniques are utilized to remove the seams in the texture maps. More details are provided in the appendix.

Synthetic Object Dataset. We assemble the dataset with 15 high-quality meshes found on the internet. We paired each object in our dataset with a detailed edit prompt to showcase our approach’s ability to perform **expressive**, **precise**, and **diverse** edits which are absent in other approaches.

Evaluation Criteria. Following Vox-E, we report auxiliary quantitative metrics on our dataset: (1) *CLIP Similarity* (CLIP_{sim}) measures the alignment of the performed 3D edits with the text descriptions, and (2) *CLIP Direction Similarity* (CLIP_{dir}) evaluates the edits with the editing directions from the input to edit results, by measuring the directional CLIP similarity between changes of text and 3D shapes, first introduced by (Gal et al. 2022).

Baselines. We compare FocalDreamer with three baselines. (1) *Fantasia3D**: as claimed in Fantasia3D, it is able to generate shapes initialized with a low-quality customized 3D mesh. In order to additionally endow it with preservation of texture from base shape, the texture field $\Gamma(p_i)$ is supervised by base texture with reconstruction loss on the base

mesh surface, as one of the baselines. (2) *Vox-E* (Sella et al. 2023): to show our superior editing within desired regions, SOTA editing work Vox-E is also compared. To the best of our knowledge, Vox-E is the only open-source method that directly performs text-guided localized edits for 3D objects. (3) *Vox-E-Global*: Vox-E also supports global editing to better align with the prompts without constraining from base shape. More details are provided in the appendix.

4.2 Qualitative Results

The qualitative comparison with 3D editing baselines is shown in Fig. 4 over our dataset. As illustrated in the figure, Fantasia3D* results in an appearance vastly different from the base mesh, even with the texture reconstruction loss, because the whole shape is re-optimized according to prompts. While Vox-E-Global occasionally produces edits that align with prompts, it suffers from subpar editing quality and noticeable outliers. Vox-E demonstrates a limited capacity to filter out undesired changes and noise based on Vox-E-Global, since it heavily relies on a keyword, such as *cape* or *glasses*. Vox-E sometimes misidentifies the focal regions, *i.e.*, placing glasses on the top of the skull. In contrast to them, our editings align perfectly with the prompts while faithfully preserving the details of base mesh, achieving precise and meaningful changes to both geometry and texture.

2D Image Editing Comparisons. We demonstrate that 2D image editing methods cannot effectively handle 3D object editing tasks, because 2D editing does not yield satisfactory view-consistent results. We sample renderings from three different viewpoints and apply SOTA image editing methods, namely Instruct Pix2Pix (IP2P) (Brooks, Holynski, and Efros 2023) and ControlNet-inpainting (ControlNet) (Zhang and Agrawala 2023). We input the same prompts in Fig. 2 for FocalDreamer, ControlNet and IP2P. As depicted in Fig. 5, the quality of editing by 2D methods drops significantly from less *canonical* views (*e.g.*, the turtle’s left view), and they severely lack view-consistency.

4.3 Quantitative Results

We perform a quantitative evaluation in Tab. 1 on our dataset. To perform a fair comparison, all metrics are calculated with renderings from the same 100 views across different methods. As illustrated in the table, FocalDreamer achieves noticeably higher CLIP_{dir}. This is attributed to its capability to accurately execute the desired editing direction, primarily due to the geometric concentration. Additionally, our editing fidelity (CLIP_{sim}) stands out as the best, stemming from the enhanced part-wise details brought by the separable framework and decoupled learning.

User Study. While CLIP mainly evaluates the matching degree of rendered views and text prompts, it fails to assess the extent to which the base shape is properly preserved. We conduct user studies with 65 participants to evaluate different methods based on user preferences across 15 scenes. We ask the participants to give a preference score (range from 1 ~ 10) in terms of prompt relevance and base shape preservation for 5 random views per scene from anonymized methods’ generation. As shown in Fig. 6, we report the distribu-

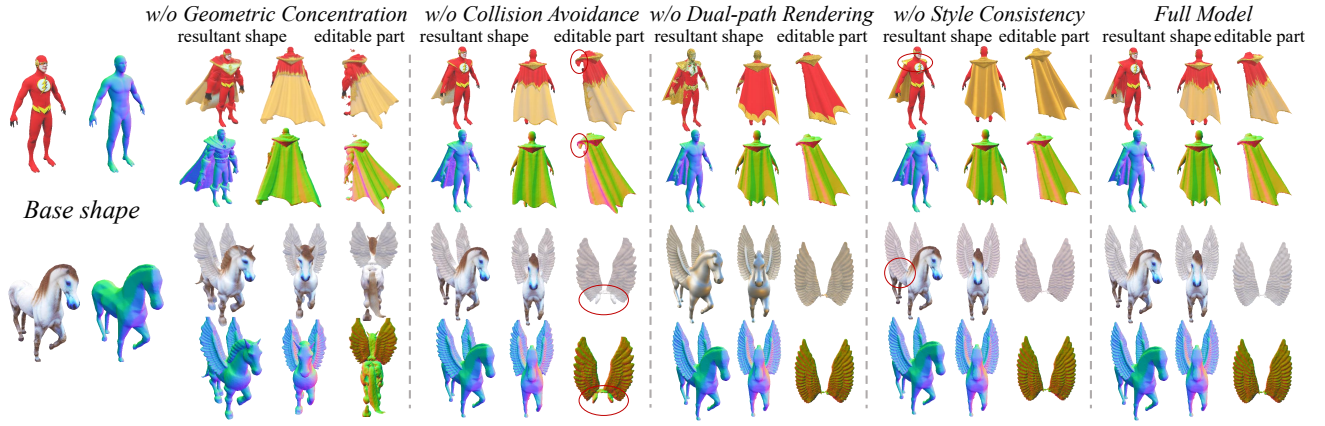


Figure 7: Ablation study. We visually illustrate the effect of each technique we propose. Please refer to Section 4.4 for details.

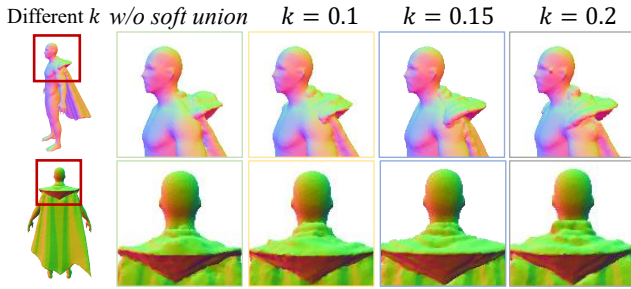


Figure 8: Geometry union sensitivity. The smoothness of the junction varies with different k in Eq. 3 and 4.



Figure 9: Progressive editing. The horse is first edited by adding two wings, then a horn is added in a subsequent edit.

tion of the scores, including the medians, means, quartiles and outliers. We find that FocalDreamer is significantly preferred over all baselines in terms of source preservation (*i.e.*, $mean = 9.14$) and prompt relevance (*i.e.*, $mean = 8.40$). The narrow interquartile range in our method also demonstrates a more stable editing effect across various scenes.

4.4 Ablation Study

We conduct the ablation study both qualitatively and quantitatively. By setting \mathcal{L}_{GF} , \mathcal{L}_{CA} and \mathcal{L}_{SC} to zero respectively, we investigated the effects of our proposed *Geometric Concentration*, *Collision Avoidance*, and *Style Consistency* strategies. To validate the *dual-path rendering*, we employ the single rendering outlined in Section 3.3. Specifically, it involves rendering the entire shape with a learnable texture Γ_e , which is initialized with the base texture Γ_b .

As illustrated in Fig. 7 and Tab. 2, \mathcal{L}_{GF} significantly constrains geometric alterations outside the focal region, resulting in localized edits. \mathcal{L}_{CA} effectively prevents undesirable geometric overlap within the base mesh, especially at the junction like the root of *wings* and *capas*. Since \mathcal{L}_{CA} predominantly affects the purity of the editable part and has minimal impact on the overall appearance, its quantitative metrics closely align with the full model. In the absence of *dual-path rendering*, the base mesh texture experiences unintended alterations due to the update of the whole texture network during appearance learning. Moreover, editing with \mathcal{L}_{SC} exhibits a harmonious overall style and nature transition in certain instances, but it is not universally required (*e.g.*, a butterfly over a tree stump). In Tab. 2, we use * to denote scenes that require \mathcal{L}_{SC} for a fair comparison.

Progressive Editing. Our method can be used as a sequential editor for users’ requirements, and progressively edits base mesh. In Fig. 9, we exhibit a two-step editing by first generating *two wings* on horse, followed by adding a *horn*.

Geometry Union Sensitivity. We also demonstrate the smoothness of the junction between the editable part and base mesh with various k (Eq. 3 and 4) in Fig. 8. It is evident that larger k leads to a more natural but pronounced transition region. We set $k = 0.15$ for a moderate transition.

5 Conclusion

In this paper, we present FocalDreamer, a text-driven framework that supports separable, precise, and consistent local editing for 3D objects. Technically, we equipped FocalDreamer with geometry union and dual-path rendering to assemble independent 3D parts, facilitating instance reuse and part-wise control. Geometric focal loss and style consistency regularization are proposed to encourage focal fusion and congruent overall appearance. Comprehensive experiments and detailed ablation studies have demonstrated our approach possesses superior local editing power through a well-conceived framework design. We hope that FocalDreamer will help pave the way for expressive, localized 3D content editing for casual artists, bringing us closer to the goal of democratizing 3D content creation for all.

Acknowledgements

This work was supported by National Science Foundation of China (U20B2072, 61976137). This work was also partly supported by SJTU Medical Engineering Cross Research Grant YG2021ZD18.

References

- Brooks, T.; Holynski, A.; and Efros, A. A. 2023. Instruct-pix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 18392–18402.
- Chen, R.; Chen, Y.; Jiao, N.; and Jia, K. 2023. Fantasia3d: Disentangling geometry and appearance for high-quality text-to-3d content creation. *arXiv preprint arXiv:2303.13873*.
- Chen, Y.; Yuan, Q.; Li, Z.; Xie, Y. L. W. W. C.; Wen, X.; and Yu, Q. 2022. Upst-nerf: Universal photorealistic style transfer of neural radiance fields for 3d scene. *arXiv preprint arXiv:2208.07059*.
- Chiang, P.-Z.; Tsai, M.-S.; Tseng, H.-Y.; Lai, W.-S.; and Chiu, W.-C. 2022. Stylizing 3d scene via implicit representation and hypernetwork. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 1475–1484.
- Dhariwal, P.; and Nichol, A. 2021. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34: 8780–8794.
- Fan, Z.; Jiang, Y.; Wang, P.; Gong, X.; Xu, D.; and Wang, Z. 2022. Unified implicit neural stylization. In *European Conference on Computer Vision*, 636–654. Springer.
- Gal, R.; Patashnik, O.; Maron, H.; Bermano, A. H.; Chechik, G.; and Cohen-Or, D. 2022. StyleGAN-NADA: CLIP-guided domain adaptation of image generators. *ACM Transactions on Graphics (TOG)*, 41(4): 1–13.
- Haque, A.; Tancik, M.; Efros, A. A.; Holynski, A.; and Kanazawa, A. 2023. Instruct-nerf2nerf: Editing 3d scenes with instructions. *arXiv preprint arXiv:2303.12789*.
- Ho, J.; Jain, A.; and Abbeel, P. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33: 6840–6851.
- Huang, Y.-H.; He, Y.; Yuan, Y.-J.; Lai, Y.-K.; and Gao, L. 2022. Stylizednerf: consistent 3d scene stylization as stylized nerf via 2d-3d mutual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 18342–18352.
- Jain, A.; Mildenhall, B.; Barron, J. T.; Abbeel, P.; and Poole, B. 2022. Zero-shot text-guided object generation with dream fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 867–876.
- Kania, K.; Yi, K. M.; Kowalski, M.; Trzciński, T.; and Tagliasacchi, A. 2022. Conerf: Controllable neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 18623–18632.
- Kobayashi, S.; Matsumoto, E.; and Sitzmann, V. 2022. Decomposing nerf for editing via feature field distillation. *Advances in Neural Information Processing Systems*, 35: 23311–23330.
- Li, J.; Niu, C.; and Xu, K. 2020. Learning part generation and assembly for structure-aware shape synthesis. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, 11362–11369.
- Li, Y.; Dou, Y.; Chen, X.; Ni, B.; Sun, Y.; Liu, Y.; and Wang, F. 2023. Generalized Deep 3D Shape Prior via Part-Discretized Diffusion Process. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 16784–16794.
- Lin, C.-H.; Gao, J.; Tang, L.; Takikawa, T.; Zeng, X.; Huang, X.; Kreis, K.; Fidler, S.; Liu, M.-Y.; and Lin, T.-Y. 2023. Magic3d: High-resolution text-to-3d content creation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 300–309.
- Lugmayr, A.; Danelljan, M.; Romero, A.; Yu, F.; Timofte, R.; and Van Gool, L. 2022. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 11461–11471.
- Metzer, G.; Richardson, E.; Patashnik, O.; Giryes, R.; and Cohen-Or, D. 2023. Latent-nerf for shape-guided generation of 3d shapes and textures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 12663–12673.
- Mikaeili, A.; Perel, O.; Cohen-Or, D.; and Mahdavi-Amiri, A. 2023. SKED: Sketch-guided Text-based 3D Editing. *arXiv preprint arXiv:2303.10735*.
- Mohammad Khalid, N.; Xie, T.; Belilovsky, E.; and Popa, T. 2022. Clip-mesh: Generating textured meshes from text using pretrained image-text models. In *SIGGRAPH Asia 2022 conference papers*, 1–8.
- Munkberg, J.; Hasselgren, J.; Shen, T.; Gao, J.; Chen, W.; Evans, A.; Müller, T.; and Fidler, S. 2022. Extracting triangular 3d models, materials, and lighting from images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 8280–8290.
- Pérez, P.; Gangnet, M.; and Blake, A. 2003. Poisson image editing. In *ACM SIGGRAPH 2003 Papers*, 313–318.
- Poole, B.; Jain, A.; Barron, J. T.; and Mildenhall, B. 2022. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*.
- Quilez, I.; and Jeremias, P. 2018. Combination SDF. <https://www.shadertoy.com/view/lt3BW2>.
- Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763. PMLR.
- Rombach, R.; Blattmann, A.; Lorenz, D.; Esser, P.; and Ommer, B. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 10684–10695.

- Sella, E.; Fiebelman, G.; Hedman, P.; and Averbuch-Elor, H. 2023. Vox-E: Text-guided Voxel Editing of 3D Objects. *arXiv preprint arXiv:2303.12048*.
- Tsalicoglou, C.; Manhardt, F.; Tonioni, A.; Niemeyer, M.; and Tombari, F. 2023. TextMesh: Generation of Realistic 3D Meshes From Text Prompts. *arXiv preprint arXiv:2304.12439*.
- Tschernezki, V.; Laina, I.; Larlus, D.; and Vedaldi, A. 2022. Neural Feature Fusion Fields: 3D distillation of self-supervised 2D image representations. In *2022 International Conference on 3D Vision (3DV)*, 443–453. IEEE.
- Xie, S.; Zhang, Z.; Lin, Z.; Hinz, T.; and Zhang, K. 2023. Smartbrush: Text and shape guided object inpainting with diffusion model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 22428–22437.
- Yang, B.; Bao, C.; Zeng, J.; Bao, H.; Zhang, Y.; Cui, Z.; and Zhang, G. 2022. Neumesh: Learning disentangled neural mesh-based implicit field for geometry and texture editing. In *European Conference on Computer Vision*, 597–614. Springer.
- Zhang, L.; and Agrawala, M. 2023. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*.
- Zhang, Y.; Huang, N.; Tang, F.; Huang, H.; Ma, C.; Dong, W.; and Xu, C. 2023. Inversion-based style transfer with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10146–10156.
- Zhuang, J.; Wang, C.; Liu, L.; Lin, L.; and Li, G. 2023. DreamEditor: Text-Driven 3D Scene Editing with Neural Fields. *arXiv preprint arXiv:2306.13455*.