

InterGSEdit: Interactive 3D Gaussian Splatting Editing with 3D Geometry-Consistent Attention Prior

Minghao Wen^{1*} Shengjie Wu^{1*} Kangkan Wang^{2†} Dong Liang^{1†}

¹MIT Key Laboratory of Pattern Analysis and Machine Intelligence,
College of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics

²The Key Lab of Intelligent Perception and Systems for
High-Dimensional Information of Ministry of Education,
School of Computer Science and Engineering, Nanjing University of Science and Technology

{lanche, wushengjie, liangdong}@nuaa.edu.cn, wangkangkan@njjust.edu.cn

Abstract

3D Gaussian Splatting based 3D editing has demonstrated impressive performance in recent years. However, the multi-view editing often exhibits significant local inconsistency, especially in areas of non-rigid deformation, which lead to local artifacts, texture blurring, or semantic variations in edited 3D scenes. We also found that the existing editing methods, which rely entirely on text prompts make the editing process a "one-shot deal", making it difficult for users to control the editing degree flexibly. In response to these challenges, we present InterGSEdit, a novel framework for high-quality 3DGS editing via interactively selecting key views with users' preferences. We propose a CLIP-based Semantic Consistency Selection (CSCS) strategy to adaptively screen a group of semantically consistent reference views for each user-selected key view. Then, the cross-attention maps derived from the reference views are used in a weighted Gaussian Splatting unprojection to construct the 3D Geometry-Consistent Attention Prior (GAP^{3D}). We project GAP^{3D} to obtain 3D-constrained attention, which are fused with 2D cross-attention via Attention Fusion Network (AFN). AFN employs an adaptive attention strategy that prioritizes 3D-constrained attention for geometric consistency during early inference, and gradually prioritizes 2D cross-attention maps in diffusion for fine-grained features during the later inference. Extensive experiments demonstrate that InterGSEdit achieves state-of-the-art performance, delivering consistent, high-fidelity 3DGS editing with improved user experience.

1. Introduction

The development of scene representation models, such as Neural Radiance Fields (NeRF) [19] and 3D Gaussian Splatting (3DGS) [15], has made scene reconstruction and rendering both efficient and practical. These advances have significantly boosted the research of 3D scene editing, and text-guided diffusion editing methods [2, 4, 6, 22, 28, 32–34] attracts increasing attention in recent years. Text-guided diffusion methods typically first edit multi-view images rendered from a 3D scene model (i.e., 3DGS) according to the target text description, and subsequently refine the 3D scene model with the modified images to generate 3D editing results. Notably, 3DGS-based scene editing [3, 4, 6, 32–34] has achieved remarkable success in rigid editing tasks such as style transfer and appearance adjustment. However, it is still challenging to edit non-rigid scenes due to more complicated and arbitrarily varying deformations.

The current mainstream approaches [4, 17, 33, 34] rely on different 3D constraint techniques, such as epipolar constraints [4], consistent 3DGS fine-tuning [33], or depth map guidance [34], to enforce cross-view consistency. Although these strategies work well for rigid editing tasks like style transfer and appearance modification, there are obvious limitations in non-rigid scene editing like faces with changing expressions. As illustrated in Fig. 1, in a "smile" editing task, the edited faces in different views may exhibit varying appearances due to non-rigid deformations in the aspects of smile intensity and tooth visibility. Refining 3DGS face with these inconsistent face images will result in blurred artifacts or distorted geometry. The primary reason for this problem is that the textual guidance in diffusion-based editing is inherently ambiguous and cannot be specified in detail to achieve consistent and precise editing. The non-rigid editing easily leads to inconsistent features across different views during the diffusion process, making it difficult to en-

*These authors contributed equally to this work.

†Corresponding author.

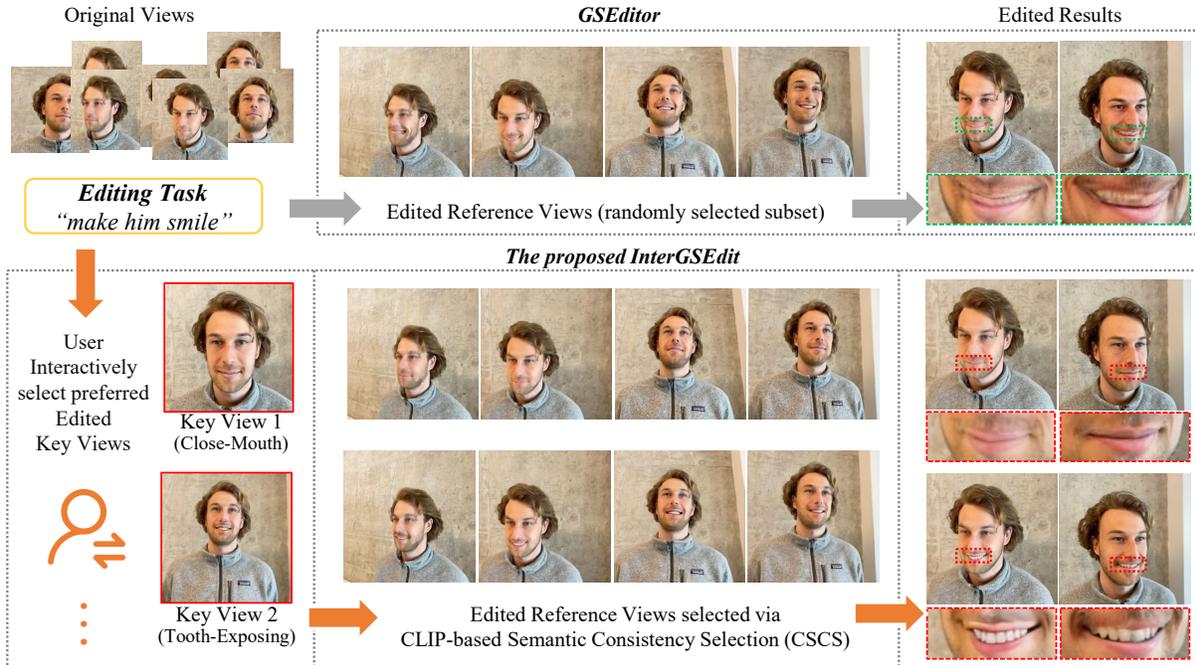


Figure 1. **Results of a “make him smile” editing task.** GSEditor [6] randomly select a subset of edited views as reference view to generate the 3DGS results, leading to blur synthesis due to multi-view inconsistency, such as tooth artifacts in this example. In contrast, our InterGSEdit framework allows the user to select the preferred key views from the edited views to guide the multi-view editing. Here, we illustrate two key views with distinct smiling characteristics (close-mouth smile and tooth-exposing smile). The edited results exhibit geometric consistency with the key views, more natural textures, and finer local details compared to GSEditor.

sure 3D geometric consistency and high-quality synthesis.

We also found that the existing 3DGS editing methods that are entirely instructed by text prompts suffer from “one-shot deal” problem in the editing process, making it difficult for users to flexibly control the editing degree. The fine-grained representation capability of the text prompts is limited by the diffusion network designation, training strategies, training sample quality, as well as users’ experience, and so on. Moreover, it is hard for the text prompts to specify fine-grained semantics due to the inherently ambiguous nature, and this linguistic uncertainty leads to variations in editing features across views, resulting in inconsistency and artifacts. In practice, achieving the 3D editing that users desire remains challenging due to the concise text prompts and the constraints of pre-trained diffusion networks.

To alleviate the above problems in 3D editing tasks, we propose a user interactive method, InterGSEdit, which builds 3D Geometry-Consistent Attention Prior (GAP^{3D}) based on a user-selected key view and ensures that the final edited results are closely consistent with the key view in both appearance and geometry. Specifically, we first perform initial editing on the rendered views, producing edited results that may exhibit discrepancies. Then, based on a key view chosen by the user, our CLIP-based Semantic Consistency Selection (CSCS) treats the key view as an anchor to select reference views with their corresponding similarity

weights that are computed from the embeddings of edited images. Finally, we construct the 3D Geometry-Consistent Attention prior (GAP^{3D}) on the 3DGS by performing a weighted unprojection of the cross-attention maps in the diffusion [6] from these selected reference views.

Then, we incorporate the constructed 3D attention prior into the same diffusion model to achieve consistent editing among different views. We project GAP^{3D} to each edited image to form a 3D-constrained attention map. To incorporate 3D geometric constraints during the denoising process, directly replacing the 2D cross-attention with the 3D-constrained attention can result in the decay of the editing details. To overcome this, we propose an adaptive cross-dimensional Attention Fusion Network (AFN) that employs a learnable gating module to adaptively fuse 2D cross-attention maps with 3D attention features. AFN ensures that during the early stages of editing, 3D-constrained attention is preferred to maintain 3D geometric consistency, while in the later stages, 2D cross attention is focused more on to preserve the fidelity of editing details, thereby achieving both structural consistency and detail recovery.

Based on InterGSEdit with these essential modules, we achieve high-quality 3D scene editing in both non-rigid tasks (e.g., facial editing), and rigid tasks (e.g., appearance modification and style transfer). Users can choose a key view to guide the editing process and obtain the final 3D

editing results that satisfy their demands on editing. Compared to existing methods, our approach is more flexible and effective to obtain view-consistent editing results and achieves high-quality synthesis and geometry. Our main contributions can be summarized in three aspects:

- We propose a user-interactive way of improving the user experience for fine-grained 3D editing and constructing 3D Geometry-Consistent Attention Prior (GAP^{3D}). After interactively selecting the key views, we employ our CLIP-based Semantic Consistency Selection (CSCS) strategy to obtain qualified reference views and then unproject their cross-attention maps to build GAP^{3D} .
- We propose an adaptive cross-dimensional Attention Fusion Network (AFN) to adaptively and dynamically fuse the GAP^{3D} with the 2D cross-attention map during the inference stages, supporting both multi-view consistency and fine-grained detail recovery.
- We propose a 3DGS editing framework, InterGSEdit, which leverages CSCS and AFN to corporately generate reliable GAP^{3D} and then achieve multi-view consistent and high-quality editing in the diffusion. Our approach demonstrates state-of-the-art performance across various scenarios on public datasets [10].

2. Related works

Image Editing with Diffusion Models. Due to the scarcity of large-scale 3D training datasets for 3D scene editing, current methods mainly utilize 2D image editing and then estimate 3D models from edited images. Among these approaches, denoising diffusion probabilistic models (DDPM) [12] are the most popular. Stable Diffusion [22] performs diffusion in a latent space, which reduces the scale of the U-Net architecture and improves the performance of diffusion models. To control the editing process, ControlNet [39] introduces control guidance such as scribbles or depth maps. GLIDE [20] leverages CLIP [21] features for image editing, while DreamBooth [24] focuses on personalized editing. Additionally, DragDiffusion [26] and Drag-a-Video [29] employ interactive mouse guidance in the editing process. Recent methods mainly focus on text-driven guidance, which is very relevant to our work. For instance, InstructPix2Pix [2] incorporates textual conditions in Stable Diffusion, and achieves high-quality text-driven editing by training on large-scale synthetic datasets. InfEdit [37] enables faithful editing for both rigid and non-rigid semantic changes using a denoising diffusion consistent model.

3D Editing. Early 3D editing methods [5, 18, 27, 40] are primarily based on traditional 3D models by synthesizing or refining meshes using text guidance. These mesh-based methods often focus on modeling a specific object rather than general scenes, resulting in a limited scope for editing. Furthermore, the editing is typically confined to modifications or optimizations on color and texture, which

constrains the applicability of such editing. With the advent of Neural Radiance Fields (NeRF) [19], subsequent approaches [9, 11, 31, 38] edit implicit scene representations, such as methods [8, 13, 35, 36] tailored for 3D character models. These methods have expanded the editing target to a relatively large scene. However, constrained by the editing capability of the guidance prior, the modifications remain confined to aspects such as color, texture, or rotation/scaling of objects. To alleviate the high computation in NeRF optimization, DreamEditor [42] utilizes view-specific masks to constrain the editing area. Also, Score Distillation Sampling [1] is introduced in [7, 14, 25] to enhance training speed and editing accuracy. More recently, diffusion-based methods edit 3D radiance fields by editing images rendered from multiple viewpoints and optimizing the underlying 3D model. With the development of the diffusion model, NeRF-based editing methods [10] incorporate diffusion models to achieve high-quality editing. The powerful ability of diffusion models enables complex 3D editing operations such as style transfer, or changing scene features such as character movements or attributes.

Multi-View Consistent 3D Gaussian Editing. Recent studies [6, 32, 41] improve accuracy, efficiency, and controllability by integrating 3DGS into 3D scene editing. Based on 3DGS, several methods are proposed to address multi-view inconsistency problem in 3D scene editing, especially for rigid editing tasks. Current methods [4, 33, 34] focus on incorporating 3D geometric constraints across multiple views to improve consistency by manipulating attention maps during the 2D editing stage. DGE [4] uses epipolar constraints to build pixel correspondences and propagate point features of the key view to the general view. VcEdit [33] reversely projects the attention maps to 3D Gaussians and then renders refined attention maps to achieve multi-view attention unification. GaussCtrl [34] minimizes the view differences by editing all views in parallel. ProEdit [3] decomposes an editing task into multiple subtasks to control the feasible output space, thereby reducing inconsistencies in the final results. These methods [4, 33] improve the consistency of multiple views to some degree, but they can lead to artifacts in non-rigid situations when the differences in 2D editing results among views are obvious. We allow users to select key views to ensure that the edited results meet their expectations, and also utilize an adaptive cross-dimensional attention fusion to maintain both view consistency and detail recovery.

3. Preliminary

3D Gaussian Splatting [15]. In 3DGS, each Gaussian, represented as G , is characterized by its mean $\mu \in \mathbb{R}^3$, covariance matrix C , associated color $c \in \mathbb{R}^3$, and opacity $\alpha \in \mathbb{R}$. The covariance matrix C can be decomposed into a scaling matrix $S \in \mathbb{R}^3$ and a rotation matrix $R \in \mathbb{R}^{3 \times 3}$, with

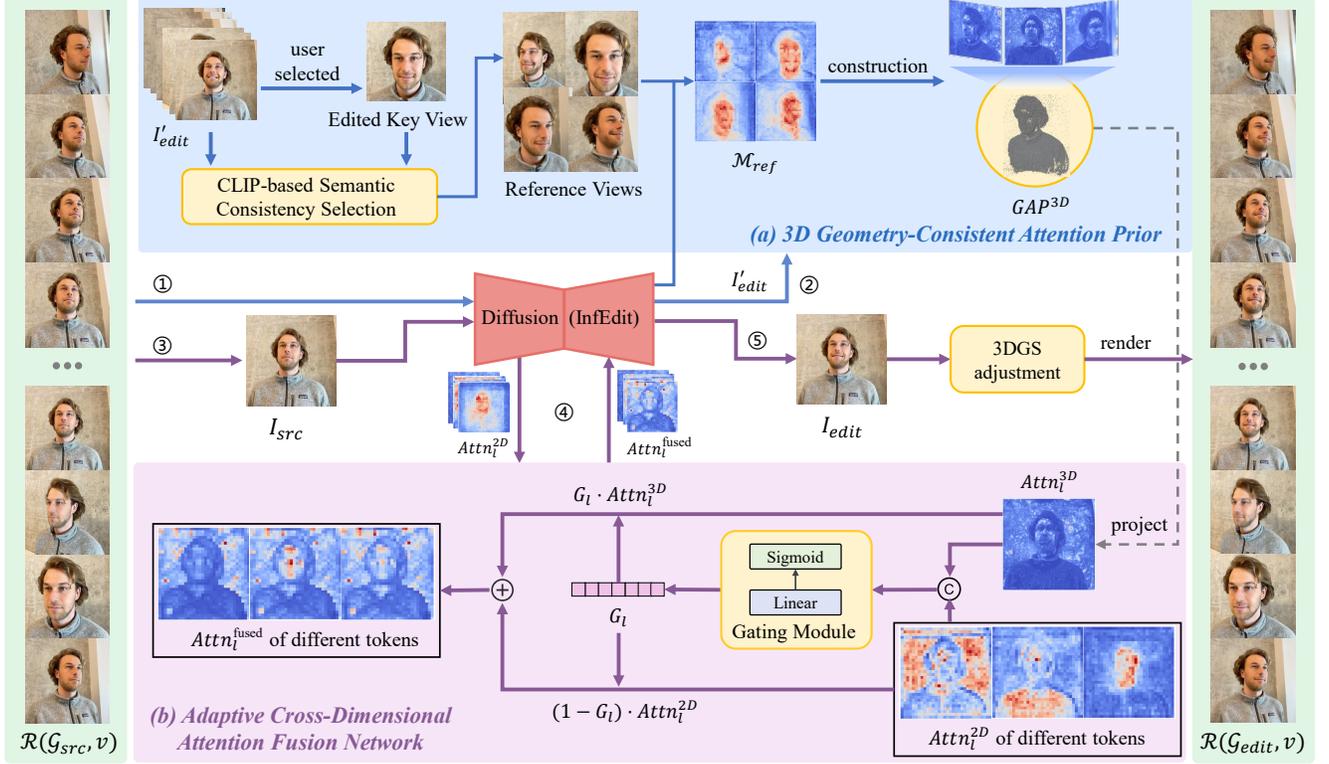


Figure 2. **Our InterGSEdit framework** mainly comprises two components. (a) 3D Geometry-Consistent Attention Prior (GAP^{3D}) Construction. User-specified key views serve to select reference views with high semantic consistency. We employ a CLIP-based Semantic Consistency Selection (CSCS) strategy to select semantically consistent reference views and then utilize their cross-attention maps to construct GAP^{3D} (Sec. 4.1). (b) Adaptive Cross-Dimensional Attention Fusion Network (AFN). We introduce a dynamic gating module to fuse 2D cross-attention with GAP^{3D} based on gating factor G_l , (Sec. 4.2), and then feed the fused results into diffusion editing.

$C = RSS^T R^T$. A Gaussian centered at μ is expressed as $G(x) = \exp(-\frac{1}{2}x^T C^{-1}x)$, where x denotes the displacement from μ to a point in 3D space. In the splatting rendering, the pixel color c is rendered by blending all sampled 3D points along the ray emitted from this pixel:

$$c = \sum_i c_i \alpha_i G(x_i) \prod_{j=1}^{i-1} (1 - \alpha_j G(x_j)). \quad (1)$$

Diffusion-based 3DGS Editing. Some methods [4, 6, 33, 34] render the 3D scene model \mathcal{G}_{src} from multiple viewpoints v to obtain 2D source images I_{src}^v . A diffusion-based editor then modifies these 2D source images into edited images I_{edit}^v according to the prompt y . By comparing the rendered images with the edited images by diffusion, an editing loss \mathcal{L}_{Edit} can be defined that measures how closely the rendered outputs of the 3D model align with the edited results guided by diffusion. When optimizing across all rendered viewpoints v , the objective function is defined as:

$$\mathcal{G}^{edit} = \arg \min_{\mathcal{G}} \sum_{v \in \mathcal{V}} \mathcal{L}_{Edit}(\mathcal{R}(\mathcal{G}, v), I_{edit}^v), \quad (2)$$

where $\mathcal{R}(\mathcal{G}, v)$ denotes differential rendering the edited \mathcal{G}

from the camera viewpoint v .

4. Methodology

Our proposed InterGSEdit aims to interactively edit a 3D scene with 3D Geometry-Consistent Attention Prior (GAP^{3D}) guided diffusion model. An illustration of our framework is shown in Fig. 2. First, the original 3D model \mathcal{G}_{src} is rendered from multiple views v to obtain source images $I_{src} = \mathcal{R}(\mathcal{G}, v)$, and I_{src} are initially edited to obtain I_{edit}^v with diffusion. Then, as described in Sec. 4.1, we introduce a CLIP-based Semantic Consistency Selection (CSCS), which leverages a user-selected key view as a semantic anchor to obtain semantically consistent reference views. These reference views, along with their associated weights, are used to construct GAP^{3D} . Subsequently, as described in Sec. 4.2, we perform 3D-constrained editing on I_{src} using a cross-dimensional Attention Fusion Network (AFN). This network employs a 3D-2D attention fusion module to inject the 3D-constrained attention, derived from GAP^{3D} projection, into the denoising process of the diffusion model. Through this diffusion editing, we obtain multi-view consistent images and the edited 3DGS.

4.1. 3D Geometry-Consistent Attention prior

4.1.1. CLIP-based Semantic Consistency Selection

The text prompts are inherently ambiguous to specify fine-grained semantics, and this linguistic uncertainty leads to variations in editing features across different views, resulting in 3D geometric inconsistency, blur outputs, and artifacts. To address this issue, we anchor the ambiguous semantics by allowing users to select an edited image as a key view, and edit the 3D scene according to the content of the key view. We propose a CLIP-based Semantic Consistency Selection (CSCS) strategy for dynamically selecting semantically consistent reference views, which quantifies the editing similarity by measuring distances in a cross-modal embedding space. Specifically, we first use the CLIP image encoder $E_{\text{CLIP}}^{\text{img}}(\cdot)$ [21] to obtain image embeddings of the edited key view $I_{\text{edit}}^{\text{key}}$ and its corresponding original image $I_{\text{src}}^{\text{key}}$, and compute their difference in editing content:

$$\Delta I_{\text{key}} = E_{\text{CLIP}}^{\text{img}}(I_{\text{edit}}^{\text{key}}) - E_{\text{CLIP}}^{\text{img}}(I_{\text{src}}^{\text{key}}), \quad (3)$$

which captures the changes in the editing direction relative to the original image. Similarly, we use the CLIP text encoder $E_{\text{CLIP}}^{\text{txt}}(\cdot)$ to obtain text embeddings of the edit prompt T_{edit} and original text T_{src} , and compute their difference as:

$$\Delta T = E_{\text{CLIP}}^{\text{txt}}(T_{\text{edit}}) - E_{\text{CLIP}}^{\text{txt}}(T_{\text{src}}), \quad (4)$$

which reflects the changes in the textual description corresponding to the editing operation. Since ΔT remains constant in denoising, it serves as a fixed reference vector representing the purpose of textual editing. We then compute the cosine similarity between ΔI_{key} and ΔT to obtain the alignment score for the key view:

$$s_{\text{key}} = D(\Delta I_{\text{key}}, \Delta T), \quad (5)$$

where D denotes the cosine distance.

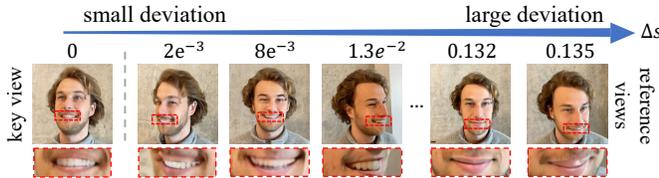


Figure 3. **Visualization of CLIP-based Semantic Consistency Selection.** The left image is the key view, whose editing direction is quantified by its alignment score s_{key} . Reference views with lower deviation value $\Delta s_v = |s_v - s_{\text{key}}|$ indicate that their editing changes are more consistent with the key view.

For a view v , we compute its alignment score as $s_v = D(\Delta I_v, \Delta T)$. Using s_{key} as an anchor, we can calculate deviation value $\Delta s_v = |s_v - s_{\text{key}}|$. The s_v exhibits editing changes that are consistent with the key view under the

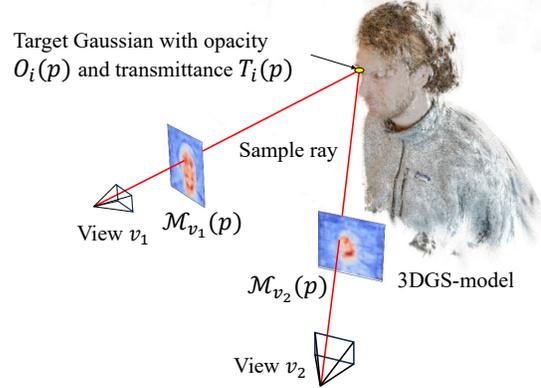


Figure 4. **Illustration of GAP^{3D} construction.** Based on the attention map \mathcal{M}_v corresponding to the view v and the associated camera pose, precise correspondences between image pixels and GS points are established. Each view has a similarity weight w_v that is calculated from Eq. (7) to construct the 3D attention prior. Here, only two views are shown for illustration.

same textual guidance if Δs_v is close to 0. We select the top- K views with minimal deviation:

$$V_{\text{ref}} = \{v_1, v_2, \dots, v_K\}, \quad (6)$$

where $\Delta s_{v_1} \leq \Delta s_{v_2} \leq \dots \leq \Delta s_{v_K}$. However, applying a hard threshold to select the top- K reference views may exclude true similar views. These views closely resemble the key view in editing features but fall just outside the top- K ranking. Such exclusions may disrupt the consistency of multi-view editing. To improve selection robustness, we introduce an adaptive weight assignment. It gives higher weights to views with greater similarity while still retaining those with slightly lower similarity if they provide useful reference information. It dynamically adjusts view weights using an exponential decay:

$$w_v = \exp(-\gamma \Delta s_v), \quad (7)$$

where w_v is the assigned weight for view v , \exp denotes the exponential function, and γ is a temperature coefficient controlling the influence of alignment deviation. If the alignment score of a candidate view s_v is close to s_{key} , its weight approaches 1, indicating a significant selection. As the alignment score deviates from s_{key} , the weight rapidly decreases, weakening its probability in the selection. This exponential decay strategy prevents useful information loss from hard-threshold filtering and retains those with slightly lower similarity, which leads to a smoother and more robust reference view selection.

4.1.2. GAP^{3D} Construction

In the following, we describe the method for constructing the 3D Geometry-Consistent Attention Prior GAP^{3D} . As illustrated in Fig. 4, given cross-attention maps \mathcal{M}_v from

viewpoint v , we first recover the 3DGS scene from reference views and then compute an attention score for each Gaussian with inverse rasterization rendering of 3DGS. Specifically, for each Gaussian point i , the attention score $GAP^{3D}(i)$ for Gaussian point i is computed as:

$$GAP^{3D}(i) = \sum_{v \in V_{ref}} \frac{w_v}{\sum_m w_m} \cdot \sum_p \mathcal{M}_v(p) \cdot O_i(p) \cdot T_i(p), \quad (8)$$

where $\mathcal{M}_v(p)$ is the attention score at pixel p in view v , $O_i(p)$ represents the opacity from Gaussian i to pixel p , and $T_i(p)$ denotes the transmittance from pixel p to Gaussian i [6]. The normalization factor $\frac{w_v}{\sum_m w_m}$ ensures that, among m reference views related to Gaussian point i , the contribution of view v is weighted according to its similarity weight. This weighting strategy on 2D attention maps establishes a 3D consistency attention prior, enhancing the consistency and quality of subsequent editing.

4.2. Cross-Dimensional Attention Fusion Network

Given a viewpoint v , we can obtain a 3D-constrained attention map formulated as $Attn^{3D} = \mathcal{R}(GAP^{3D}, v)$ by projecting 3D attention prior to the 2D image space. By injecting 3D-constrained attention into UNet [23] of the diffusion model, we can integrate 3D geometry constraints across multiple attention layers. Here, we design a cross-dimensional Attention Fusion Network (AFN) that fuses 3D-constrained attention $Attn^{3D}$ with the 2D cross-attention obtained in the editing process. This method can facilitate the stability and consistency of editing features throughout the denoising process. Specifically, after the cross-attention map $Attn_l^{2D}$ is generated at layer l within the diffusion model [2, 17, 22], we project 3D attention prior to a 3D-constrained attention map $Attn_l^{3D}$ and fuse with the cross-attention map $Attn_l^{2D}$.

Directly replacing 3D-constrained attention with 2D cross-attention may lead to feature degradation during editing diffusion. To address this problem, AFN adaptively adjusts the influence of 3D-constrained attention at different stages with a dynamic gated fusion mechanism. Specifically, we introduce a gating factor G_l to fuse the features of $Attn_l^{3D}$ and $Attn_l^{2D}$ dynamically. To ensure that the model emphasizes 3D geometric consistency during the initial inference phase and gradually focuses on editing details later, we design a linearly decaying dynamic bias term given by $\gamma(t) = \alpha(1 - \frac{t}{T})$, where α is a constant, t denotes the current iteration, and T is the total number of iterations.

$$G_l = \sigma\left(W_l \cdot [Attn_l^{2D}; Attn_l^{3D}] + \gamma(t)\right), \quad (9)$$

$$Attn_l^{fused} = G_l \cdot Attn_l^{3D} + (1 - G_l) \cdot Attn_l^{2D}, \quad (10)$$

where W_l is a learnable weight matrix for the attention fusion at different layers, and σ is the sigmoid function mapping the gating factor G_l into the interval $(0, 1)$. We further

incorporate a KL divergence constraint to enhance training stability and precisely regulate the attention fusion process between $Attn^{3D}$ and $Attn^{2D}$. The total optimization objective is defined as:

$$\mathcal{L}_{total} = \lambda_{2D} \mathcal{L}_{Edit} + \lambda_{3D} \mathcal{L}_{KL}(Attn^{3D} \parallel Attn^{2D}), \quad (11)$$

where \mathcal{L}_{Edit} represents the editing loss as the diffusion objective in Eq. (2), and λ_{2D} and λ_{3D} are hyperparameters for two loss terms. \mathcal{L}_{KL} enforces distributional consistency between the 2D cross-attention $Attn^{2D}$ and 3D-constrained attention $Attn^{3D}$. This constraint forces the $Attn^{2D}$ distribution to converge to the $Attn^{3D}$ distribution, thus ensuring that the final generation results are geometry-consistent.

To progressively adjust the attention fusion during inference, we adopt a dynamic regulation strategy that simultaneously optimizes the Attention Fusion Network (AFN) and the scene 3DGS in the denoising process. Specifically, the weights W_l and a Gating Module are learned to obtain the gating factor G_l for the attention fusion, while a KL divergence constraint is used to regulate the fusion process. In the early stages, the model is encouraged to prioritize geometric consistency by setting a larger weight λ_{3D} of the KL loss, which forces AFN to adjust $Attn^{2D}$ to better align with $Attn^{3D}$. As denoising process proceeds and the geometric information stabilizes, by gradually reducing λ_{3D} , the model is enforced to adjust G_l so that $Attn^{2D}$ contributes more to the denoising guided by text prompts, ensuring fine editing details are restored. Consequently, our approach achieves the dual objectives of enforcing structural consistency during early training and progressively emphasizing detail recovery.

5. Experiments

5.1. Experimental Settings

Implementation Details. We implement our framework based on ThreeStudio [16]. We utilize the diffusion model of InfEdit [37] as the foundation for diffusion-based editing and employ the latent consistency model [17] from HuggingFace. To demonstrate our method’s capability in 3D model editing, we conduct experiments on various scenes from IN2N datasets [10]. We represent 3D scenes using 3DGS [15] and leverage segmentation pipeline and semantic tracking of GSEditor [6] for local editing. Each editing is performed on a set of 20 random views, and the optimization of 3DGS is conducted for 800 to 1200 iterations depending on the complexity of the scene. All experiments are conducted on a NVIDIA RTX 4090 GPU.

Evaluation Metrics. We evaluate our method both qualitatively and quantitatively. The quantitative evaluation is based on the following metrics: CLIP Similarity [21], CLIP Text-Image Direction Similarity (CTIDS) [3], and CLIP

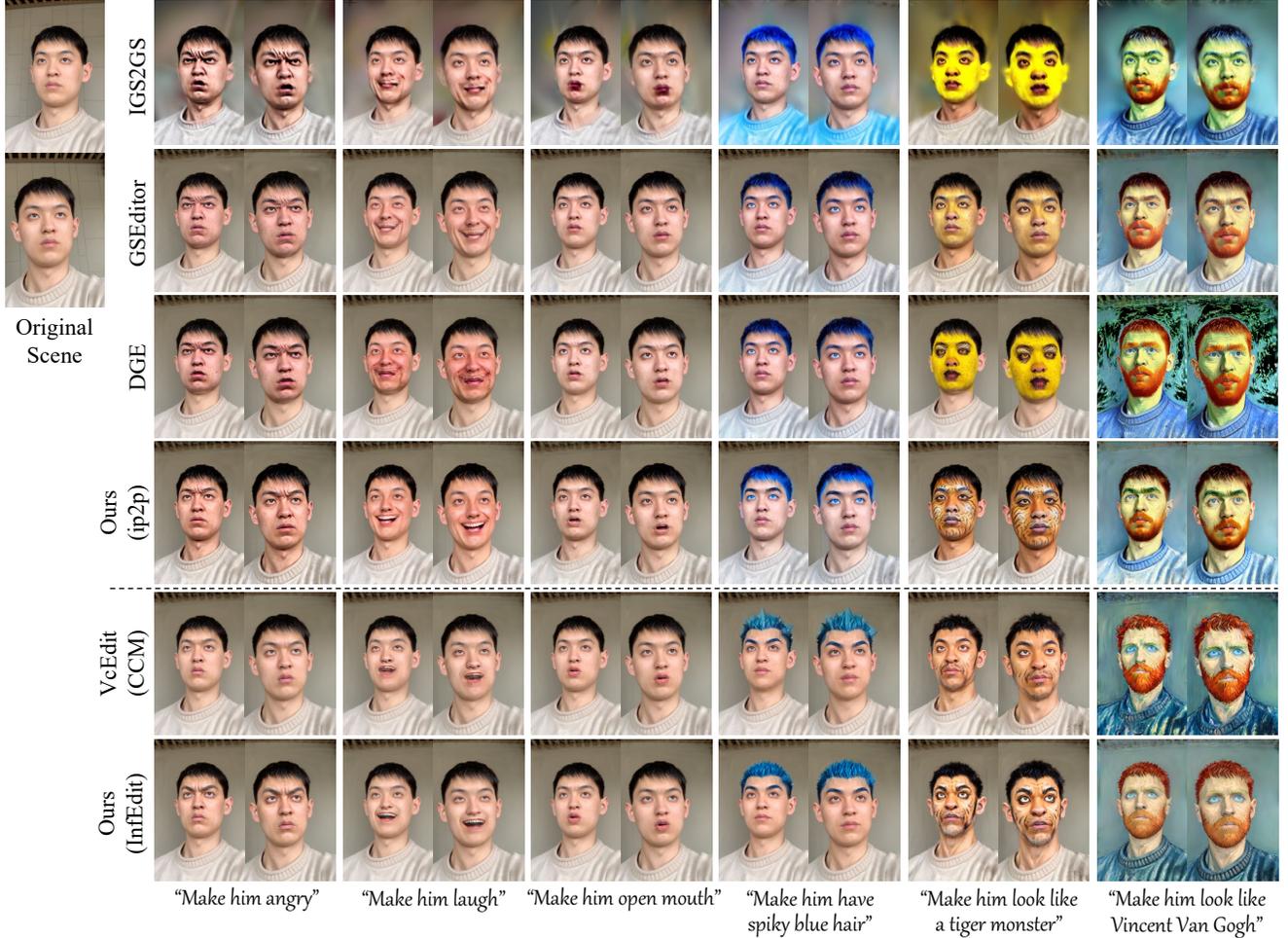


Figure 5. **Qualitative comparison** with IGS2GS [30], GSEditor [6], DGE [4] and VcEdit [33] with CCM module. Our InterGSEdit framework achieves high-quality editing for both non-rigid and rigid tasks, demonstrating strong fidelity to textual instructions, and precise geometry preservation. In contrast to other methods that produce edited results with tooth artifacts and unnatural modifications, our approach significantly reduces such artifacts, enhances texture quality, and produces more natural editing.

Direction Consistency (CDC) [10]. Specifically, we randomly sample 20 camera views from the 3DGS training dataset for multi-view editing and fine-tune the 3DGS to obtain the final model with the edited images. For quantitative evaluation, we render multi-view images from the resulting 3D model using all available camera poses and assess the performance across the three evaluation metrics. The CLIP similarity score is computed as the cosine similarity between the CLIP-encoded edited image embedding $E_{\text{CLIP}}^{\text{img}}(I_{\text{edit}})$ and the target text embedding $E_{\text{CLIP}}^{\text{txt}}(T_{\text{edit}})$. CTIDS [3] is calculated as the cosine similarity between the textual embedding difference $\Delta T = E_{\text{CLIP}}^{\text{txt}}(T_{\text{edit}}) - E_{\text{CLIP}}^{\text{txt}}(T_{\text{src}})$ and the corresponding image embedding difference $\Delta I = E_{\text{CLIP}}^{\text{img}}(I_{\text{edit}}) - E_{\text{CLIP}}^{\text{img}}(I_{\text{src}})$. Additionally, following IN2N [10], CDC is used to measure directional consistency in editing results.

5.2. Qualitative Comparisons

To demonstrate the superiority of our InterGSEdit, we compare it with several baseline methods: IGS2GS [30], GSEditor [6], DGE [4] and VcEdit [33] with CCM module. For IGS2GS, GSEditor, and DGE, we implement them with their released codes and models. Since the codes of VcEdit [33] are not released publicly, we fulfill its cross-attention consistency module (CCM) and perform diffusion editing of InfEdit [37] to compare with their CCM module on consistent editing. Note that IGS2GS, GSEditor, and DGE adopt the InstructPix2Pix (ip2p) [2] as their editing backbone, while VcEdit employs the InfEdit [37]. To ensure a fair comparison and demonstrate the robustness of our approach across different diffusion frameworks, we conduct experiments using both ip2p and InfEdit for our method. Fig. 5 presents qualitative comparisons in the Fangzhou scene.

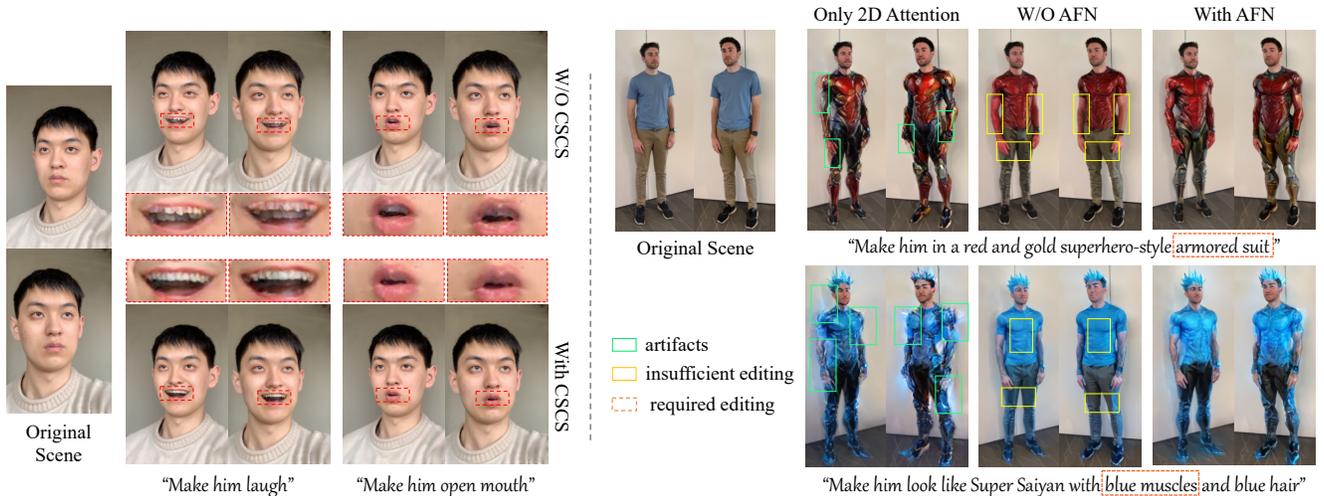


Figure 6. **Ablation studies of our CLIP-based Semantic Consistency Selection (CSCS) and Attention Fusion Network (AFN).** In the version “W/O CSCS”, non-rigid editing results, such as facial expression editing, exhibit serious tooth artifacts due to inconsistent features among different editing views. And the version “W/O AFN” ignore the editing features of clothing, causing the humans are not edited properly and the texture details are not modified to match the textual description.

Method	CLIP Similarity \uparrow	CTIDS \uparrow	CDC \uparrow
IGS2GS [30]	0.2166	0.1328	0.8071
GSEditor [6]	0.2169	0.1372	0.8099
DGE [4]	0.2098	0.1097	0.8118
VcEdit [33]	0.2195	0.1195	0.8043
Ours (ip2p)	0.2265	0.1416	0.8424
Ours (InfEdit)	0.2285	0.1531	0.8347

Table 1. **Quantitative comparison** of different methods. Higher values indicate better performance. Notably, our InterGSEdit outperforms competing approaches in terms of CLIP similarity, CTIDS, and CDC metrics.

In the nonrigid cases, our InterGSEdit view-consistently produces more natural and realistic expressions, as demonstrated in the “make him angry” and “make him laugh” tasks. In contrast, existing methods tend to introduce undesirable artifacts, such as tooth artifacts, attributed to varying tooth features across multi-view editing results. For appearance and style transfer tasks, our method also outperforms other methods in editing quality. For example, when editing a man’s hair to be blue and spiky hair, our method produces more natural colors and hair textures that accurately reflect the “spiky” semantic characteristics.

5.3. Quantitative Comparisons

As shown in Tab. 1, our method achieves the highest scores across all the metrics, demonstrating superior text-image alignment and multi-view consistency. Specifically, our method achieves a CLIP Similarity score of 0.2285, outperforming GSEditor (0.2169) and DGE (0.2098). Similarly,

	CLIP Similarity \uparrow	CTIDS \uparrow	CDC \uparrow
W/O CSCS	0.2034	0.1301	0.7647
With CSCS	0.2090	0.1475	0.8218
Only 2D Attention	0.2252	0.2835	0.8616
W/O AFN	0.2206	0.1243	0.8429
With AFN	0.2403	0.2738	0.8802

Table 2. **Ablation study of our CLIP-based Semantic Consistency Selection (CSCS) and Attention Fusion Network (AFN).** The results demonstrate the effectiveness of CSCS and AFN in achieving high-quality and multi-view consistent editing results.

our method achieves a CTIDS score of 0.1531, surpassing GSEditor (0.1372) and DGE (0.1097), indicating that our approach better captures the semantic intent of the editing prompt. Additionally, our method attains the highest CDC score (0.8347), confirming improved 3D-aware editing consistency across different views.

5.4. Ablation Study

We evaluated the crucial role of CSCS and AFN within InterGSEdit. In the version without CSCS, a 3D attention prior is constructed by averaging unprojection from all views. As shown in Fig. 6, “W/O CSCS” results in obvious tooth artifacts due to geometry inconsistencies across different edited views. In the version without AFN, 3D-constrained attention is directly substituted for 2D cross-attention maps to generate the final outputs. And in the “Only 2D Attention” version, 2D cross-attention features from the diffusion model remain unaltered during the editing process. As shown in Fig. 6, compared to the “With

AFN” version, the version “W/O AFN” exhibits some regions that are not fully edited due to the lack of original editing features, and the version “Only 2D Attention” shows clear artifacts and blurred regions without the constraint of 3D geometry prior. Furthermore, without AFN, the much lower accuracy of the CTIDS metric indicates that the details are not recovered and the editing is not consistent to textual instructions. Notably, the “Only 2D Attention” version shows a slightly higher CTIDS score than ours because this version utilizes full 2D cross-attention, leading to the most extensive editing. In contrast, our AFN fuses 3D constraint attention with 2D cross attention, making its CTIDS score theoretically a bit lower or equal. Quantitatively, the improvements in CTIDS and CDC metrics indicate that our framework using CSCS and AFN not only aligns with the desired editing semantics, but also achieves multi-view consistency, as demonstrated in Tab. 2.

6. Conclusion

In this work, we introduce InterGSEdit, an innovative 3D editing approach based on 3D Gaussian Splatting (3DGS) that enhances the quality and consistency of 3D editing. In our framework, anchored with a user-specified key view, CLIP-based Semantic Consistency Selection is used to construct a 3D Geometry-Consistent Attention Prior from the selected reference views. In the subsequent multi-view editing process, 3D Geometry-Consistent Attention Prior is projected into 3D-constrained attention maps for all views, and then dynamically fused with the 2D cross-attention maps via the Attention Fusion Network. This fusing attention effectively guides the generation of editing diffusion to preserve geometric consistency and capture fine-grained details, producing coherent and realistic editing results.

7. Acknowledgments

This work is supported in part by the National Natural Science Foundation of China under Grant 62272229 and 62472224, the Natural Science Foundation of Jiangsu Province under Grant BK20222012.

References

- [1] Thiemo Alldieck, Nikos Kolotouros, and Cristian Sminchisescu. Score distillation sampling with learned manifold corrective. In *European Conference on Computer Vision*, pages 1–18. Springer, 2024. 3
- [2] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18392–18402, 2023. 1, 3, 6, 7
- [3] Jun-Kun Chen and Yu-Xiong Wang. Proedit: Simple progression is all you need for high-quality 3d scene editing. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 1, 3, 6, 7
- [4] Minghao Chen, Iro Laina, and Andrea Vedaldi. Dge: Direct gaussian 3d editing by consistent multi-view editing. In *European Conference on Computer Vision*, pages 74–92. Springer, 2024. 1, 3, 4, 7, 8
- [5] Yongwei Chen, Rui Chen, Jiabao Lei, Yabin Zhang, and Kui Jia. Tango: Text-driven photorealistic and robust 3d stylization via lighting decomposition. *Advances in Neural Information Processing Systems*, 35:30923–30936, 2022. 3
- [6] Yiwen Chen, Zilong Chen, Chi Zhang, Feng Wang, Xiaofeng Yang, Yikai Wang, Zhongang Cai, Lei Yang, Huaping Liu, and Guosheng Lin. Gaussianeditor: Swift and controllable 3d editing with gaussian splatting. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 21476–21485, 2024. 1, 2, 3, 4, 6, 7, 8
- [7] Xinhua Cheng, Tianyu Yang, Jianan Wang, Yu Li, Lei Zhang, Jian Zhang, and Li Yuan. Progressive3d: Progressively local editing for text-to-3d content creation with complex semantic prompts. In *The Twelfth International Conference on Learning Representations*, 2024. 3
- [8] Ruikai Cui, Xibin Song, Weixuan Sun, Senbo Wang, Weizhe Liu, Shenzhou Chen, Taizhang Shang, YANG LI, Nick Barnes, Hongdong Li, et al. Lam3d: Large image-point clouds alignment model for 3d reconstruction from single image. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 3
- [9] Ori Gordon, Omri Avrahami, and Dani Lischinski. Blended-nerf: Zero-shot object generation and blending in existing neural radiance fields. In *2023 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, page 2933–2943. IEEE, 2023. 3
- [10] Ayaan Haque, Matthew Tancik, Alexei A Efros, Aleksander Holynski, and Angjoo Kanazawa. Instruct-nerf2nerf: Editing 3d scenes with instructions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19740–19750, 2023. 3, 6, 7
- [11] Runze He, Shaofei Huang, Xuecheng Nie, Tianrui Hui, Luoqi Liu, Jiao Dai, Jizhong Han, Guanbin Li, and Si Liu. Customize your nerf: Adaptive source driven 3d scene editing via local-global iterative training. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6966–6975, 2024. 3
- [12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 3
- [13] Fangzhou Hong, Mingyuan Zhang, Liang Pan, Zhongang Cai, Lei Yang, and Ziwei Liu. Avatarclip: zero-shot text-driven generation and animation of 3d avatars. *ACM Transactions on Graphics (TOG)*, 41(4):1–19, 2022. 3
- [14] Hiromichi Kamata, Yuiko Sakuma, Akio Hayakawa, Masato Ishii, and Takuya Narihira. Instruct 3d-to-3d: Text instruction guided 3d-to-3d conversion. *arXiv preprint arXiv:2303.15780*, 2023. 3
- [15] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkuehler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (TOG)*, 42(4):1–14, 2023. 1, 3, 6

- [16] Ying-Tian Liu, Yuan-Chen Guo, Vikram Voleti, Ruizhi Shao, Chia-Hao Chen, Guan Luo, Zixin Zou, Chen Wang, Christian Laforte, Yan-Pei Cao, et al. Threestudio: A modular framework for diffusion-guided 3d generation. *ICCV*, 2023. 6
- [17] Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthesizing high-resolution images with few-step inference. *arXiv preprint arXiv:2310.04378*, 2023. 1, 6
- [18] Oscar Michel, Roi Bar-On, Richard Liu, Sagie Benaim, and Rana Hanocka. Text2mesh: Text-driven neural stylization for meshes. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13492–13502, 2022. 3
- [19] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 1, 3
- [20] Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *International Conference on Machine Learning*, pages 16784–16804. PMLR, 2022. 3
- [21] 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. 3, 5, 6
- [22] 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, 3, 6
- [23] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015. 6
- [24] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22500–22510, 2023. 3
- [25] Etai Sella, Gal Fiebelman, Peter Hedman, and Hadar Averbuch-Elor. Vox-e: Text-guided voxel editing of 3d objects. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 430–440, 2023. 3
- [26] Yujun Shi, Chuhui Xue, Jun Hao Liew, Jiachun Pan, Hanshu Yan, Wenqing Zhang, Vincent YF Tan, and Song Bai. Dragdiffusion: Harnessing diffusion models for interactive point-based image editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8839–8849, 2024. 3
- [27] Yawar Siddiqui, Antonio Alliegro, Alexey Artemov, Tatiana Tommasi, Daniele Sirigatti, Vladislav Rosov, Angela Dai, and Matthias Nießner. Meshgpt: Generating triangle meshes with decoder-only transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19615–19625, 2024. 3
- [28] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative gaussian splatting for efficient 3d content creation. In *The Twelfth International Conference on Learning Representations*, 2024. 1
- [29] Yao Teng, Enze Xie, Yue Wu, Haoyu Han, Zhenguo Li, and Xihui Liu. Drag-a-video: Non-rigid video editing with point-based interaction. *arXiv preprint arXiv:2312.02936*, 2023. 3
- [30] Cyrus Vachha and Ayaan Haque. Instruct-gs2gs: Editing 3d gaussian splats with instructions, 2024. 7, 8
- [31] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Clip-nerf: Text-and-image driven manipulation of neural radiance fields. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3835–3844, 2022. 3
- [32] Junjie Wang, Jiemin Fang, Xiaopeng Zhang, Lingxi Xie, and Qi Tian. Gaussianeditor: Editing 3d gaussians delicately with text instructions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 20902–20911, 2024. 1, 3
- [33] Yuxuan Wang, Xuanyu Yi, Zike Wu, Na Zhao, Long Chen, and Hanwang Zhang. View-consistent 3d editing with gaussian splatting. In *European Conference on Computer Vision*, pages 404–420. Springer, 2024. 1, 3, 4, 7, 8
- [34] Jing Wu, Jia-Wang Bian, Xinghui Li, Guangrun Wang, Ian Reid, Philip Torr, and Victor Adrian Prisacariu. Gaussctrl: Multi-view consistent text-driven 3d gaussian splatting editing. In *European Conference on Computer Vision*, pages 55–71. Springer, 2024. 1, 3, 4
- [35] Kailu Wu, Fangfu Liu, Zhihan Cai, Runjie Yan, Hanyang Wang, Yating Hu, Yueqi Duan, and Kaisheng Ma. Unique3d: High-quality and efficient 3d mesh generation from a single image. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 3
- [36] Ziyi Wu, Yulia Rubanova, Rishabh Kabra, Drew Hudson, Igor Gilitschenski, Yusuf Aytar, Sjoerd van Steenkiste, Kelsey Allen, and Thomas Kipf. Neural assets: 3d-aware multi-object scene synthesis with image diffusion models. *Advances in Neural Information Processing Systems*, 37: 76289–76318, 2025. 3
- [37] Sihan Xu, Yidong Huang, Jiayi Pan, Ziqiao Ma, and Joyce Chai. Inversion-free image editing with language-guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9452–9461, 2024. 3, 6, 7
- [38] Yu-Jie Yuan, Yang-Tian Sun, Yu-Kun Lai, Yuewen Ma, Rongfei Jia, and Lin Gao. Nerf-editing: geometry editing of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18353–18364, 2022. 3
- [39] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In

Proceedings of the IEEE/CVF international conference on computer vision, pages 3836–3847, 2023. [3](#)

- [40] Qihang Zhang, Chaoyang Wang, Aliaksandr Siarohin, Peiye Zhuang, Yinghao Xu, Ceyuan Yang, Dahua Lin, Bolei Zhou, Sergey Tulyakov, and Hsin-Ying Lee. Towards text-guided 3d scene composition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6829–6838, 2024. [3](#)
- [41] Shijie Zhou, Haoran Chang, Sicheng Jiang, Zhiwen Fan, Zehao Zhu, Dejie Xu, Pradyumna Chari, Suyu You, Zhangyang Wang, and Achuta Kadambi. Feature 3dgs: Supercharging 3d gaussian splatting to enable distilled feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21676–21685, 2024. [3](#)
- [42] Jingyu Zhuang, Chen Wang, Liang Lin, Lingjie Liu, and Guanbin Li. Dreameditor: Text-driven 3d scene editing with neural fields. In *SIGGRAPH Asia 2023 Conference Papers*, pages 1–10, 2023. [3](#)