# GSEditPro: 3D Gaussian Splatting Editing with Attention-based Progressive Localization

Y. Sun[†1] , R. Tian[†1] , X. Han[†1] , X. Liu[2] , Y. Zhang[‡1] and K. Xu[2]

[1]State Key Laboratory for Novel Software Technology of Nanjing University, China
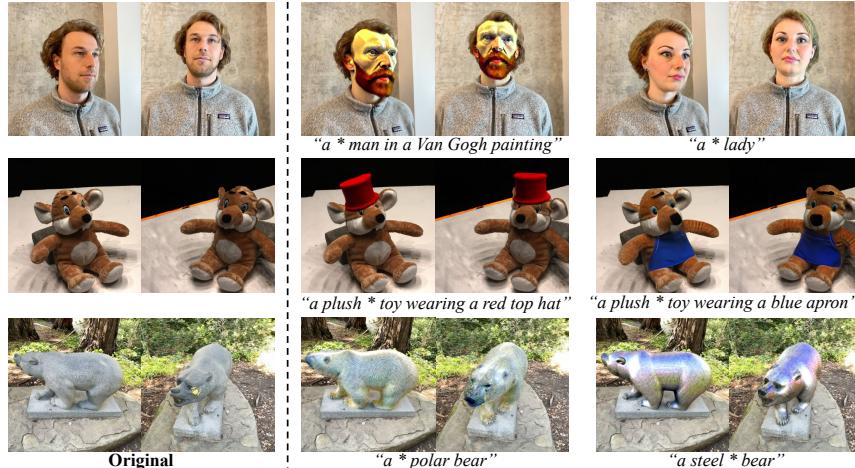[2]National University of Defense Technology, China

**Figure 1:** *Results of GSEditPro. GSEditPro enables users to conduct high-quality editing in various scenes using text prompts only.*

**Abstract**
*With the emergence of large-scale Text-to-Image(T2I) models and implicit 3D representations like Neural Radiance Fields (NeRF), many text-driven generative editing methods based on NeRF have appeared. However, the implicit encoding of geometric and textural information poses challenges in accurately locating and controlling objects during editing. Recently, significant advancements have been made in the editing methods of 3D Gaussian Splatting, a real-time rendering technology that relies on explicit representation. However, these methods still suffer from issues including inaccurate localization and limited manipulation over editing. To tackle these challenges, we propose GSEditPro, a novel 3D scene editing framework which allows users to perform various creative and precise editing using text prompts only. Leveraging the explicit nature of the 3D Gaussian distribution, we introduce an attention-based progressive localization module to add semantic labels to each Gaussian during rendering. This enables precise localization on editing areas by classifying Gaussians based on their relevance to the editing prompts derived from cross-attention layers of the T2I model. Furthermore, we present an innovative editing optimization method based on 3D Gaussian Splatting, obtaining stable and refined editing results through the guidance of Score Distillation Sampling and pseudo ground truth. We prove the efficacy of our method through extensive experiments.*

**CCS Concepts**
• *Computing methodologies* → *Rendering; Point-based models; Computer vision representations;*

---

† Equal contributions to this work.

‡ Corresponding author.

## 1. Introduction

In the rapidly evolving field of computer graphics, developing user-friendly methods for 3D generation and editing is crucial, as these methods can be widely applied in domains such as virtual reality and digital gaming.

In the field of 3D generation, text-based model generation technology [PJBM22, LGT*22, MRP*23, WLW*24] has made significant progress due to the success of large-scale Text-to-Image (T2I) models [SCS*22, RDN*22, YXK*22]. These models demonstrate remarkable creativity and significantly reduce the cost of model generation, gaining increasing attention. However, they often lack editing abilities, and even slight text prompt variations may lead to different output results. In addition to 3D generation tasks, editing existing 3D models is also crucial, which enables efficient and precise modifications to 3D models, thereby increasing the flexibility and adaptability of the existing models for various applications.

In recent years, the emergence of implicit 3D representation Neural Radiance Fields (NeRF) [MST*21] has made significant progress in scene reconstruction and novel view synthesis. The high-fidelity rendering ability of NeRF and its significant scalability provide excellent support for subsequent work. Consequently, most text-driven 3D editing techniques [HTE*23, MPS*23, WCH*22] have been designed based on NeRF for quite some time. However, editing neural radiance fields is difficult due to its implicit encoding of shape and texture information in high-dimensional neural network features. Thus, accurate locating and direct modification during the editing process are challenging, hindering the obtainment of precise and high-quality editing results, thereby impeding their practical applications. A pioneering work recently emerging in the field of 3d reconstruction is 3D Gaussian Splatting (3D-GS) [KKLD23], a real-time rendering technology based on explicit representation. The explicit nature of 3D-GS gives it a significant advantage in editing tasks. Each 3D Gaussian distribution exists independently, allowing for editing 3D scenes easily by directly manipulating the 3D Gaussians required for editing constraints. Recently, some editing methods [CCZ*23, FWZ*23, YDYK23, ZKC*24] based on 3D-GS have emerged. However, they still encounter various issues such as inaccurate locating or requiring users to manually locate editing areas in some cases [ZKC*24, CCZ*23], difficulty in ensuring consistency of non-editing areas before and after editing [YDYK23, CCZ*23], inability to perform object insertion operations effectively [CCZ*23, FWZ*23], and failure to guarantee consistency between different viewpoints after editing [FWZ*23].

To overcome these issues, we propose a novel text-driven editing framework based on 3D-GS called GSEditPro, which enables users to perform 3D editing intuitively and precisely using text prompts. Our framework achieves this through two key designs: (1) Attention-based editing area localization in 3D: We leverage the explicit representation advantage of 3D-GS to classify Gaussians based on their relevance to the attention maps derived from cross-attention layers of T2I models, assigning semantic labels to each Gaussian, thereby obtaining accurate 3D editing mask areas. (2) Guidance for a detailed Optimization from DreamBooth and pseudo-GT images: We create an optimization process that balances generative capability and detail preservation. It uses simple text prompts to effectively perform 3D scene editing by conducting score distillation sampling within the 3D editing mask area, thus ensuring high-quality editing. Additionally, we maintain the details of the scenes by constructing pseudo-GT images to ensure consistency of the irrelevant regions using pixel-level guidance.

We conducted experiments using the proposed method in various synthetic and real-world 3D scenes. The experiments demonstrate that our editing method can achieve precise editing both on object changing and object insertion with irrelevant areas naturally preserved after editing. Furthermore, since editing is accomplished through simple text prompts, our method is highly user-friendly, showcasing significant practical application potential. Qualitative and quantitative comparisons also indicate that our method outperforms previous methods in terms of editing accuracy, visual fidelity, and user satisfaction.

Our contributions can be summarized as follows:

1. We propose GSEditPro, a novel 3D editing method that enables users to perform various creative and precise editing operations using only text prompts. This approach is more convenient than previous Gaussian editing methods, which require additional user input as prior. Our extensive experiments demonstrate that our framework still offers advantages in both qualitative and quantitative metrics with user-friendly interactions.

2. We design a method to add semantic labels to Gaussians using the cross-attention mechanism of the T2I model and the explicit representation advantage of 3D-GS. Our special attention-based localization module assists in achieving more accurate 3D editing area localization and more effective editing control.

3. We present how to preserve details with pixel-level guidance, which creates a pseudo-GT image using our localization module to minimize unnecessary modifications and guide 3D Gaussian rendering for more detailed results.

## 2. Related Work

### 2.1. Text-guided Image Generation and Editing

Today, numerous methods have attained impressive outcomes in the realm of text-driven image generation and editing. T2I diffusion models [HSC*22, RBL*22, SCS*22] based on large-scale image-text data demonstrate diverse and high-quality image generation capabilities that align well with text prompts. However, these models do not offer the ability to modify the generated images. Prompt-to-Prompt [HMT*22] utilizes images provided by users to generate images based on text prompts while editing them simultaneously, providing an intuitive editing method. DreamBooth [RLJ*23] adjusts diffusion models using an L2 reconstruction loss and proposes a preservation loss to avoid overfitting. It obtains extensive updated model parameters, providing the capability to generate high-quality images and perform editing. Additionally, some methods [WDR*24, ZRA23]incorporate spatial conditions that control object positions during the generation process. This enhances the model's ability to handle various input conditions, allowing for finer control over the generated image results. These methods have achieved excellent results in the 2D domain, but extending them directly from 2D to 3D is not straightforward.

## 2.2. Text-to-3D

With the development of T2I generation models, interest in the Text-to-3D domain is continuously increasing. However, directly applying diffusion models in the 3D domain is a challenging task since it requires keeping the consistency of different views. Dream-Field [JMB*22] employs the image-text embedding model CLIP [RKH*21] to guide the optimization of NeRF [MST*21], successfully generating 3D shapes from text prompts. DreamFusion [PJBM22] first proposed the score distillation sampling (SDS) loss, which obtains priors from a pre-trained T2I model and optimizes the neural radiance field during training. Based on DreamFusion, a series of works [LGT*22, MRP*23, RKP*23] adopted a similar optimization process and achieved better generation results by refining the process and employing different SDS guidance methods. Furthermore, recent works [ZRX*24, DYL*24, LSZ*24]utilize 3D-GS [KKLD23] as their 3D representation, enabling rapid generation of 3D models based on text prompts. However, the generation results of these methods are easily influenced by the effectiveness of text prompts, and they are limited to generating 3D models, unable to edit existing 3D scenes.

## 2.3. Text-guided 3D Editing

Text-guided neural radiance field editing has gained significant attention as a new research area. EditNeRF [LZZ*21] pioneered this field, offering the ability to edit both the shape and color of NeRF based on implicit encoding. Subsequently, some methods [MPS*23, WCH*22] began combining NeRF [MST*21] and diffusion models. For instance, Instruct-NeRF2NeRF [HTE*23] utilized a text-based diffusion model [BHE23]) to modify rendered images based on the user's instructions and gradually modify the neural radiance field, achieving excellent editing effects. However, due to the implicit representation of NeRF, these editing methods lack precise control over editing regions. Therefore, previous methods [ZWL*23, SFHAE23] adopted explicit representation methods such as grids and voxels to improve the quality of local editing. However, these methods have not obtained satisfactory editing results for real-world scenes. Concurrent work ConsistentDreamer [CBM*24] adds 3D-consistent structured noise to rendered multi-view images, and applies self-supervised consistency training using consistency-warped images, generating 3D consistently edited images from 2D diffusion models.

The introduction of 3D-GS [KKLD23] presents an opportunity to overcome this limitation. Its explicit representation enhances control over editing regions. GaussianEditor [CCZ*23] utilizes 3D-GS as scene representation and employs semantic tracking technology to track the cloning and splitting of Gaussians, resulting in more accurate scene editing results. Another GaussianEditor [FWZ*23] leverages large language models to extract Regions of Interest (RoI) from text prompts and converts them to the image space using segmentation models. Then it trains them using the loss proposed in SA3D [CZF*23] to elevate RoI to the 3D scene. GaussianGrouping [YDYK23] adds identity encoding to Gaussians, classifying them and providing the ability to modify 3D objects. However, the editing results of these methods lack consistency across different viewpoints, and they fail to add objects to the scene or request additional input from users, which results

in significant inconvenience. Concurrent works aim to address this problem. DGE [CLV24] injects features from selected key views into the diffusion network through correspondence matching with epipolar constraints, enabling direct editing views of a 3D model with a text-based image editor. From the perspective of image editing, VCEdit [WYW*24] intergrates two view-consistent modules into the Gaussian editing framework. TIGER [XCC*24] introduces Coherent Score Distillation, which combines a 2D image editing diffusion model with a multi-view diffusion model for score distillation, resulting in multi-view consistent outcomes.

## 3. Our Method

GSEditPro is a novel editing framework designed to edit a pre-trained 3D-GS scene according to the provided text prompt. It alters the geometry and appearance of the objects of interest within the original scene while ensuring that the 3D content unrelated to the prompt remains unchanged.

The overall framework of GSEditPro, as illustrated in Figure 2, consists of two main stages. Firstly, we design an attention-based localization module that employs a T2I model and the cross-attention mechanism to locate the editing region in the 3D space using the keywords in the text prompt, as elaborated in Section 3.2. Secondly, building upon 3D Gaussians, we implement scene editing leveraging Score Distillation Sampling (SDS) loss with Dreambooth [RLJ*23], as detailed in Sections 3.3. We employ the attention-based localization module at the pixel level to further preserve the unrelated areas during the editing process, as depicted in Section 3.4. By integrating optimization and progressive localization, our method achieves precise and detailed local editing.

### 3.1. Preliminary

**3D Gaussian Splatting.** 3D Gaussian Splatting (3D-GS) [KKLD23] is an explicit representation method for 3D scenes, utilizing a set of anisotropic 3D Gaussians to represent the scene, denoted as $G = g_1, g_2, ..., g_N$, where $g_i = \{\mu, \Sigma, c, \alpha\}$, $i \in 1, ..., N$. Among them, $\mu$ denotes the center position of the Gaussian, $\Sigma$ represents the 3D covariance matrix, $c$ is the RGB color represented by spherical harmonic coefficients, and $\alpha$ denotes opacity. 3D-GS employs a differentiable splatting rendering method, enabling high-quality real-time rendering. The splatting rendering process can be formulated as:

$$C = \sum_{i \in \mathcal{N}} c_i \sigma_i \prod_{j=1}^{i-1}(1 - \sigma_j) \qquad (1)$$

where $\mathcal{N}$ represents the number of Gaussians contributing to the ray, $\sigma_i = \alpha_i e^{-\frac{1}{2}(x_i)^T \Sigma^{-1}(x_i)}$ denotes the influence of the Gaussian on the image pixel, and $x_i$ is the distance between the pixel and the center of the $i$-th Gaussian.

**SDS Loss.** DreamFusion [PJBM22] first introduced the Score Distillation Sampling (SDS) loss, which guides the generation of 3D models by extracting prior knowledge from the T2I diffusion model. It first adds noise at level $t$ to a randomly rendered view $I$ to obtain $I_t$. Then it uses a pre-trained diffusion model $\Phi$ to predict the added noise under the condition of $I_t$ and the text prompt $y$. The
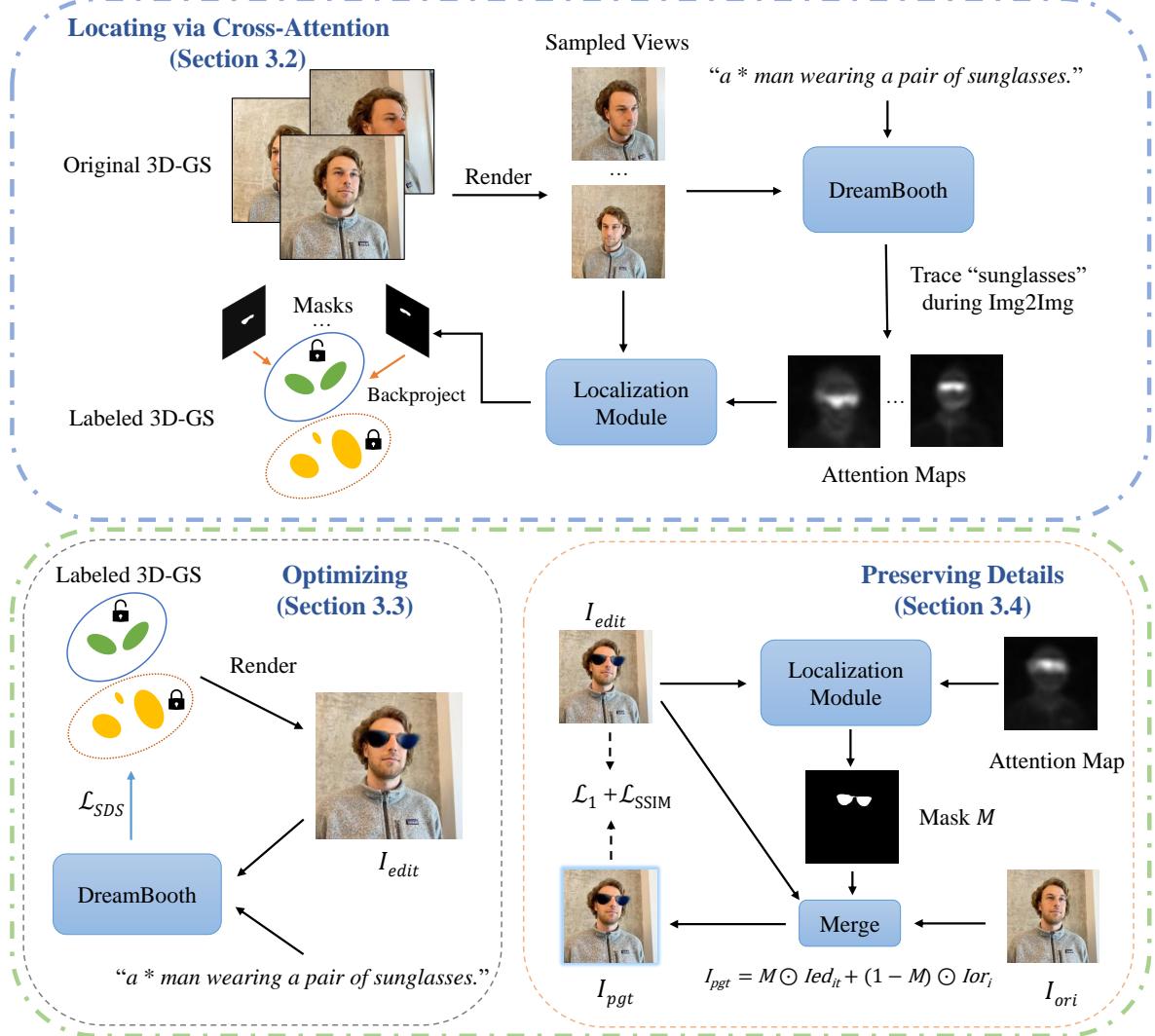
**Figure 2:** *Method Overview. GSEditPro edits scenes represented by 3D Gaussian Splatting using text prompt only. The key process of our method contains three parts: 1) Locating Gaussian editing regions via cross-attention, which assigns semantic labels to each Gaussian and determines whether the gradients can be propagated between them (Section 3.2); 2) Optimizing editing regions using DreamBooth, which uses $\mathcal{L}_{SDS}$ as the guidance to optimize Gaussian parameters iteratively (Section 3.3); 3) Preserving details with pixel-level guidance, which creates a pseudo-GT image to guide 3D Gaussian rendering for more detailed results (Section 3.4).*

SDS loss is calculated as the gradient for each pixel as follows:

$$\nabla_\theta \mathcal{L}_{SDS}(\Phi, I = g(\theta)) = \mathbb{E}_{\epsilon,t}[w(t)(\epsilon_\Phi(I_t : y, t) - \epsilon)\frac{\partial I}{\partial \theta}] \quad (2)$$

where $w(t)$ is a weighting function based on the noise level $t$, $\theta$ represents the parameters of the neural radiance field, and $g$ denotes the rendering process function. During training, gradients are back-propagated to $\theta$, guiding the rendering results of the neural radiance field to be more similar to the images generated by the T2I diffusion model based on text prompts.

### 3.2. Locating Gaussian Editing Regions via Cross-Attention

Initially, we locate the region of interest for modification within 3D space based on text prompts, laying the foundation for the editing framework. Previous 3D editing approaches [LLF*23, KMS22] have commonly relied on static 2D masks for determining the editing region, yet such methods suffer from limitations in ensuring consistency of the edited scene across views. With dynamic changes occurring in the 3D representation during training, these masks can become inaccurate or even ineffective. Additionally, some works [XH22, YSL*22] have utilized 3D masks for locating. Still, they need manual operation to get the 3D masks in some

cases, which is conducted with complex rules and may result in the imprecision of the located region.

The cross-attention layers inside the T2I diffusion model can capture the relationship between the generated image and each word [HMT*22]. Similarly, during our editing process, we need to manipulate the target objects within the 3D Gaussians under the control of text prompts through the T2I diffusion model. Therefore, we propose an attention-based localization module that utilizes the 2D probability maps generated by the cross-attention layers as masks for each view, determining which regions need editing in 2D views. These 2D maps are then processed as point prompts for the large-scale segmentation model Segment Anything Model(SAM) [KMR*23] to obtain a more precise mask for the target region. After that, we backproject the 2D masks into 3D space and mark the Gaussians that need editing, enabling precise localization of the editing region explicitly in Gaussians before training.

Concretely, we sample rendering output in various views using COLMAP [SF16] cameras and fine-tune the Stable Diffusion [RBL*22] using DreamBooth [RLJ*23]. DreamBooth is a method that fine-tunes the large-scale text-to-image (T2I) model around a specific target subject, denoted as "*" or other symbols, to ensure its ability to generate images similar to the input data. To strengthen the generating stability and ability of the fine-tuned diffusion model, we set the class prompt as the target editing prompt. The preservation loss of DreamBooth will encourage the diffusion model to treat this special class as the default generating style, which increases the accuracy of attention maps as well.

Furthermore, we collect the attention maps of the target words during the Img2ImgPipe of DreamBooth, which generates several images based on our editing prompt. These maps from cross-attention layers represent the rough or possible position of the editing area depending on whether it exists in the original scene, which means our method can have a reasonable localization of incorporation editing with the prior of the Diffusion model. Our localization module will decide how to locate the region according to the existence of the target object. Suppose the text prompt for editing is about adding new objects to the scene. In that case, the localization module chooses the filtered attention maps directly as the preliminary results, and the threshold is set as 0.5 in our experiments. Considering maps lack precision, preliminary results are clustered using the clustering algorithm DBSCAN [EKSX96] to filter out outliers further to get the final 2D masks. When editing existing objects in the scene, our module first tries to use Language-based SAM [KMR*23] to segment them in sampled views. However, the results of the SAM based on the language prompt differ from the views which will result in bad results of the masks. And it always fails to segment the part of the target editing objects. So we will improve the results when they have a small overlap over the attention maps. The traced maps are filtered and then clustered as mentioned before. The localization module chooses points of the processed maps as point prompts for the SAM, with the top 5 points selected based on the highest attention map values as positive ones, while the negative point prompts are chosen based on the lowest 3 values. After that SAM will segment a precise mask of the target for each view. Masks are back-projected during the differentiable rendering process similar to GaussianEditor [CCZ*23] and we only allow

gradients to propagate within the labeled Gaussians whose weights of back-projection bigger than the threshold. Finally, our method finishes Locating Gaussian editing regions explicitly and assigns the Gaussians their binary labels in 3D.

### 3.3. Optimizing Editing Regions using DreamBooth

After locating the editing regions, we propose an optimization scheme for 3D Gaussian editing. To achieve text-based 3D editing, we introduce the diffusion model in the optimization stage. After training on our target dataset, DreamBooth [RLJ*23] possesses sufficient generation ability to guide the training of 3D Gaussians. We utilize the SDS loss proposed by DreamFusion [PJBM22] as the guiding loss function. After obtaining the prompt for editing and the images rendered from random views during training, they are collectively used as inputs to compute $\mathcal{L}_{SDS}$ in DreamBooth. This loss is then employed during the back-propagation process to guide the cloning and splitting of the Gaussians, as well as the changes in their parameters. The computation can be formulated as follows:

$$\mathcal{L}_{SDS} = w(t)(\epsilon_{\Phi}(x_t, t, T) - \epsilon)^2 \tag{3}$$

where $w(t)$ is the weight of SDS decided by timestep $t$, $\epsilon_{\Phi}$ is the denoiser of the diffusion model to compute the noise which will be removed, $x_t$ is the embedding vector of the noised image, $T$ is the text prompt input. To conveniently control all losses in the method with weight, $\mathcal{L}_{SDS}$ adopts the squared error between real noise and predicted noise.

In Section 3.2, we have already finished locating editing regions and only allowed the gradients to be backpropagated between the Gaussians to be edited. Therefore, during each training iteration, $\mathcal{L}_{SDS}$ serves as a 2D guidance to optimize Gaussian parameters iteratively. This process matches the rendering results with the text guidance of the editing, obtaining desired editing results after sufficient training.

Inspired by GaussianEditor [CCZ*23], the parameters of the GS model are supposed to be constrained according to its existing generations, which will prevent it from being exposed to the randomness of the loss function. And it is represented as:

$$L_{anchor}^{P} = \sum_{i=0}^{n} \lambda_i (P_i - \hat{P}_i)^2 \tag{4}$$

where $P$ denotes the property of Gaussians including the position $x$, the scaling $s$, the rotation $q$, the transparent $\alpha$ and the color $c$. And $n$ is the maximum generation now. $\hat{P}$ represents the saved anchor state. $\lambda_i$ refers to the strength of the loss applied, which will grow with the increase of the number of terms.

### 3.4. Preserving Details with Pixel-Level Guidance

The strong fluidity of Gaussians makes it easy to cause editing beyond the desired region. We propose a method to maintain the accuracy of the editing region at the pixel level. As mentioned in Section 3.2, our approach has already locked the Gaussian regions to be edited after locating, with gradients only passing through the target Gaussians. However, due to the strong fluidity of Gaussians,
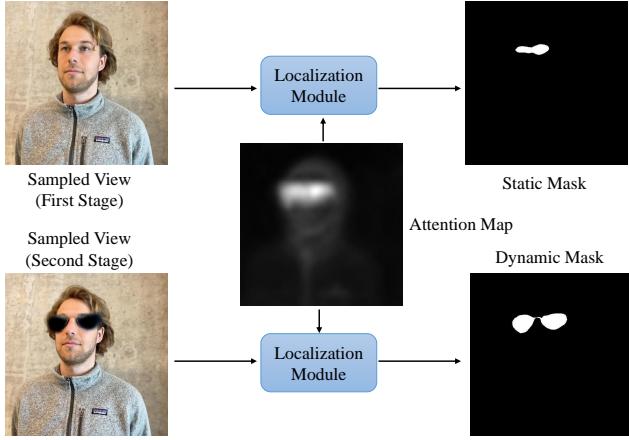
**Figure 3:** *The visual difference of the Mask between the two stages. The static Mask is a guessed sketch from the DreamBooth, and the dynamic mask is located in real-time after 3D-GS capable of rendering decent edited results.*

there may still be changes outside the editing region in rendered results. We create a pseudo-GT image for each rendered image during training to ensure consistency between the editing results and the original images. The pseudo-GT image is generated by combining the modified parts in the current rendered result with the unedited parts of the initial rendered result as shown in the bottom right of Figure 2. Note that in our overview the pseudo-GT image is almost the same as the rendered image because they are similar to each other indeed when our method converges. We then use $\mathcal{L}_1$ and $\mathcal{L}_{D-SSIM}$ losses to constrain the similarity between the current rendered result and the pseudo-GT image, ensuring that after back-propagation, the overall result shifts towards greater consistency.

In order to get a pseudo-GT image to guide 3D Gaussian rendering for desired editing results, we need an appropriate mask to separate the editing region from others accurately. We divide the generation of this mask into two stages as shown in Figure 3, with the main difference lying in the method of selecting masks.

In the first stage, as shown in Section 3.2, a mask suitable for locating the editing region can be obtained for each rendered view through the localization module. Therefore, during the initial 2000 iterations of rendering, we utilize this static mask to construct a coarse pseudo-GT image.

In the second stage, when the editing results of the Gaussian rendering have a roughly formed shape, we reuse the localization module introduced in Section 3.2 to locate a dynamic mask for generating a more reliable pseudo-GT image, which changes dynamically during the training optimization. The calculation method for obtaining the pseudo-GT image through masking is as follows:

$$I_{pgt} = M \odot I_{edit} + (1 - M) \odot I_{ori} \qquad (5)$$

where $I_{pgt}$ is the pseudo-GT image, $M$ is the obtained mask, $I_{edit}$ is the edited image, and $I_{ori}$ is the original image. After obtaining $I_{pgt}$,

we calculate the pixel-wise loss with the edited image as follows:

$$\mathcal{L}_{Preservation} = \lambda_{l1}\mathcal{L}_1(I_{edit}, I_{pgt}) + \lambda_{ssim}\mathcal{L}_{D-SSIM}(I_{edit}, I_{pgt}) \qquad (6)$$

where $\lambda_{l1}$ and $\lambda_{ssim}$ represent the weighting values assigned to $\mathcal{L}_1$ and $\mathcal{L}_{D-SSIM}$ respectively. Experiment results indicate that our pixel-wise editing consistency preservation method achieves optimal results at this stage. For specific details, please refer to the ablation experiments.

Thus, the loss function $\mathcal{L}$ of our method can be expressed as:

$$\mathcal{L} = \lambda_{sds}\mathcal{L}_{SDS} + \mathcal{L}_{Preservation} + \sum_{P \in \{x,s,q,\alpha,c\}} \lambda_P L_{anchor}^P \qquad (7)$$

where $\lambda_{sds}$ is the weighting value assigned by $\mathcal{L}_{SDS}$ defined in Equation 3 and $\lambda_P$ is the weighting value for different anchor loss of the Gaussian parameters $x, s, q, \alpha, c$.

## 4. Experiments

### 4.1. Experimental Setup

**Dataset.** We use scenes including human faces, indoor settings, and complex outdoor environments, providing a comprehensive evaluation of our method's efficacy. For outdoor environments, we use the Mip-NeRF360 [BMV*22] dataset; for faces and indoor settings, we use datasets provided in GaussianEditor [CCZ*23] and DreamEditor [ZWL*23]. We employ the training method from 3D-GS [KKLD23] and camera viewpoints selected from COLMAP [SF16] to train the original Gaussians. For each editing task, we use a text prompt about the scene as input and select a keyword to fine-tune the diffusion model.

**Baseline.** We compare our approach with three baselines. Earlier text-based editing works often rely on NeRF [MST*21], so we select two representative text-based neural radiance field editing methods: Instruct-NeRF2NeRF [HTE*23](I-N2N) and DreamEditor [ZWL*23]. I-N2N utilizes Instruct-Pix2Pix [BHE23] to update rendered multi-view images based on specific text instructions. DreamEditor employs a mesh-based representation and incorporates DreamBooth [RLJ*23] to support text-based editing. Additionally, for the latest editing works based on 3D-GS [KKLD23], we choose the state-of-the-art GaussianEditor [CCZ*23] for comparison. Similar to I-N2N, GaussianEditor utilizes Instruct-Pix2Pix to iteratively optimize parameters of Gaussians to guide 3D-GS in completing editing tasks.

**Evaluation Criteria.** Following I-N2N and GaussianEditor, we use CLIP [RKH*21] text-image Directional Similarity ($CLIP_{dir}$). $CLIP_{dir}$ assesses the alignment between changes in text pairs and corresponding changes in image pairs. Detailed definitions will be presented in supplementary material. To ensure a fair comparison, we standardize prompts and instructions into the same format and generate the original text using the BLIP-2 [LLSH23] model.

Considering that the quality of editing is closely related to human perception, which cannot be fully quantified by $CLIP_{dir}$, we have also conducted user studies. We presented multi-view images of the original and edited scenes to participants and gathered their preferences. During the survey, we randomized the order of results so that users were unaware of which result belonged to which method.

**Figure 4:** *More multi-view results of our method in different scenes. The left column is the original view and the other three columns are the multi-view editing results. Our method is capable of conducting various kinds of editing in different scenes. Different from most previous methods, our method succeeds in altering the geometry and appearance of objects detailedly.*



**Figure 5:** *Comparisons on adding objects to the given scene. Results of the GaussianEditor shown above are generated without manually adjusting the estimated depth during training. Our method gives attention to proper locations and generates satisfactory results. In contrast, GaussianEditor [CCZ\*23] incorrectly positions Gaussians, and I-N2N [HTE\*23] fails to edit as instructed.*

**Implementation Details.** For baseline comparisons, we mostly follow the recommended settings in their papers, except when extending the training iterations is necessary. For GaussianEditor and I-N2N, we make comparisons with a total of 14 editing tasks on 4 scenes. We collect 48 questionnaires for the user study. For DreamEditor, considering that they have not published their pre-process method, we will conduct an additional comparison, using a subset of their released preprocessed datasets and checkpoints, with a total of 7 editing tasks on 3 scenes. We collect 56 questionnaires for the user study.

### 4.2. Qualitative Results

We provide qualitative results of our method in Figure 1 and Figure 4. Results show that our method can properly locate editing regions and edit various scenes, which succeeds in altering the geometry and appearance of objects based on the text prompts. Since

GaussianEditor [CCZ\*23] implements different pipelines for object insertion and other editing tasks, we will first make comparisons on addition with GaussianEditor and I-N2N [HTE\*23]. As shown in Figure 5, we present results about adding sunglasses to the plush toy and a birthday hat to the man. GaussianEditor first selects a single view, repaints it, and makes monocular depth estimation, which cannot guarantee a precise depth map. So it tends to insert objects into the wrong locations. Besides, some unintended details can be observed(e.g. the face under the hat in the middle column of the first row). Although dynamically adjusting estimated depth can alleviate the first problem, the second one remains unsolved. To be supplemented, it is noticed that the details of the inserted object are not fine enough. I-N2N reveals its issue that it misunderstands prompts and generates confusing results, due to the limited ability of Instruct-Pix2Pix. In contrast, our method neither generates undesirable results nor mistakenly inserts objects into unexpected regions, leading to satisfactory and clean results.
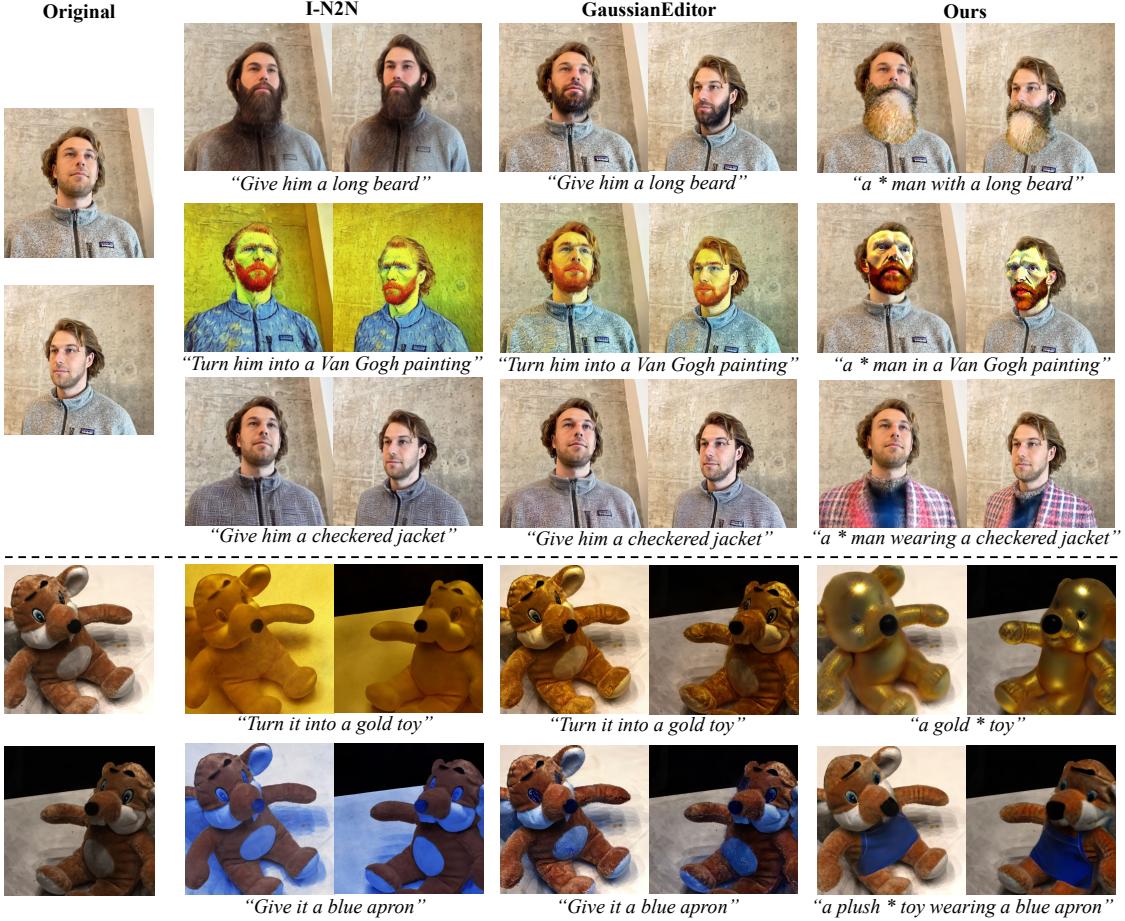
**Figure 6:** *Comparisons with GaussianEditor [CCZ\*23] and I-N2N [HTE\*23]. We show superior ability in making successful edits and controlling editing regions, without leaking noisy Gaussians to backgrounds and irrelative regions.*

As presented in Figure 6, we make comparisons on editing tasks. As mentioned before, I-N2N sometimes fails to handle complex instructions properly and cannot generate results corresponding to editing prompts (in the third row and the last row). Besides, it can be observed that editing regions spread all over the scene (in the second row and the fourth row) since I-N2N utilizes Instruct-Pix2Pix to fulfill dataset updates, which makes modifications to the whole picture. The second column shows the editing results of GaussianEditor, which demonstrates its limitations. In the first row, we try to generate a long beard for him, but the result doesn't meet our expectations. A possible reason is that its semantic labels cannot fully correspond with the prompt. However, our optimization guided by DreamBooth [RLJ\*23] does not need to focus on this issue. Our pseudo-GT images prevent editing results from changing backgrounds (in the second row), while carefully designed guidance promises sufficient changes according to the prompt (in the fourth row). Moreover, GaussianEditor inherits the shortcomings of I-N2N, as shown in the third and the last row, failing to edit correctly when instructions are complex.

Noticing that the editing pipelines in both GaussianEditor and I-N2N utilize Instruct-Pix2Pix, we also selected DreamEditor [ZWL\*23] for another comparison. Results are shown in Figure 7. DreamEditor also fully leverages the capabilities of DreamBooth [RLJ\*23], achieving relatively high-quality edits on scenes. However, we can generate a beard with better shape (in the first row) and a hat with more details (in the second row). We also control the edit region, ensuring the dog does not appear noisy (in the last row). We believe that the key to our better results lies in the fact that driving Gaussians to the target region is easier than manipulating a mesh-based NeRF.

### 4.3. Quantitative Results

We present the $CLIP_{dir}$ score and votes of users in Table 1 and Table 2. In Table 1, we compare our method with GaussianEditor [CCZ\*23] and I-N2N [HTE\*23]. The results clearly indicate that our method achieves significantly higher $CLIP_{dir}$ scores and wider preference, suggesting better alignment between editing texts and results in our method.

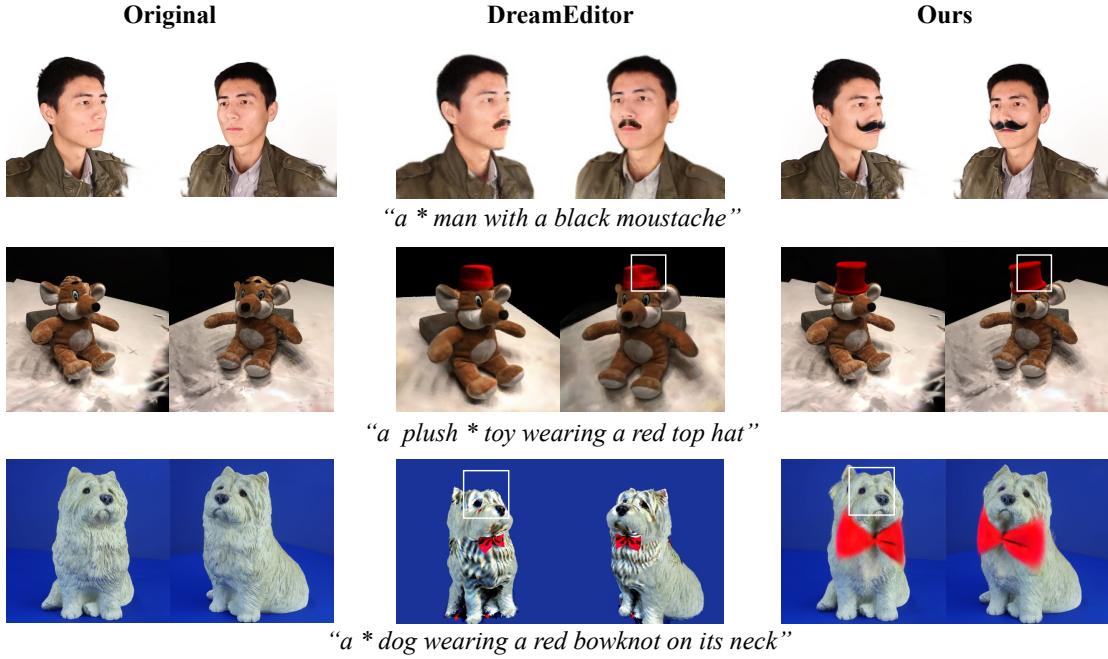In Table 2, we compare with DreamEditor [ZWL\*23]. The re-

| Original | DreamEditor | Ours |
|----------|-------------|------|



*"a * man with a black moustache"*

*"a  plush * toy wearing a red top hat"*

*"a * dog wearing a red bowknot on its neck"*

**Figure 7:** *Comparisons with DreamEditor [ZWL*23]. Our method adds a better-shaped mustache on the man, and a more detailed red top hat on the plush toy. For the dog, we make a more precise editing. Please be aware that the fragmented clothing in the first scene and the rough edges in the second scene result from the limited viewpoints used to train the original Gaussian models. This is an issue with the original Gaussians, but not with our editing method.*

**Table 1:** *Quantitative comparison with GaussianEditor [CCZ*23] and I-N2N [HTE*23] in CLIP Directional Similarity metrics and user study evaluations.*

| Method | $CLIP_{dir}\uparrow$ | Vote↑ |
|--------|---------------------|-------|
| I-N2N | 0.1542 | 23.1% |
| GaussianEditor | 0.1467 | 22.5% |
| Ours | **0.1883** | **54.4%** |

sults indicate that our method achieves significantly higher $CLIP_{dir}$ scores, showing that the shapes and textures generated by our method are more consistent with the editing text prompts and with better quality. In conclusion, we found that the evaluation results further demonstrate that GSEditPro achieves higher user satisfaction in various scenarios, which also proves that the users prefer to more fantastic edit more than just texture edit.

**Table 2:** *Quantitative comparison with DreamEditor [ZWL*23] in CLIP Directional Similarity metrics and user study evaluations.*

| Method | $CLIP_{dir}\uparrow$ | Vote↑ |
|--------|---------------------|-------|
| DreamEditor | 0.1911 | 38.1% |
| Ours | **0.2021** | **61.9%** |

### 4.4. Ablation Study

**Effectiveness of our attention-based localization of editing areas.** To demonstrate the effectiveness of the precise localization of editing regions proposed in this paper, we conduct two experiments. (1) Without localization: We omit the localization step in Section 3.2 and directly optimize all Gaussians with the full steps afterward. (2) Our method: By adding semantic labels to Gaussian distributions, we precisely determine the editing region based on attention maps and only optimize the selected editing region during the optimization process. As shown in Figure 8, the method without localization inadvertently alters irrelevant areas in the scene(e.g. the plants near the stump), disrupting the consistency of non-editing areas. In contrast, with precise editing region localization, our method ensures that changes occur only in the regions of interest.

**Effectiveness of the pixel-level guidance using pseudo-GT images**. Our detailed framework is accomplished by conducting score distillation sampling within the 3D editing mask region first and then further optimizing the quality with pseudo-GT images. Therefore, in ablation experiments, we design the following experiments. (1)Without pixel-level guidance: We remove the pixel-level guidance of pseudo-GT images and use the SDS guidance only. (2) Our method: Adding pseudo-GT images in training for further precise editing. As shown in Figure 9, it can be observed that in the results without pixel-level guidance, although we have located the editing regions in 3D as a constraint, which tries to enable the reduction of unnecessary modifications in irrelevant areas, the background is
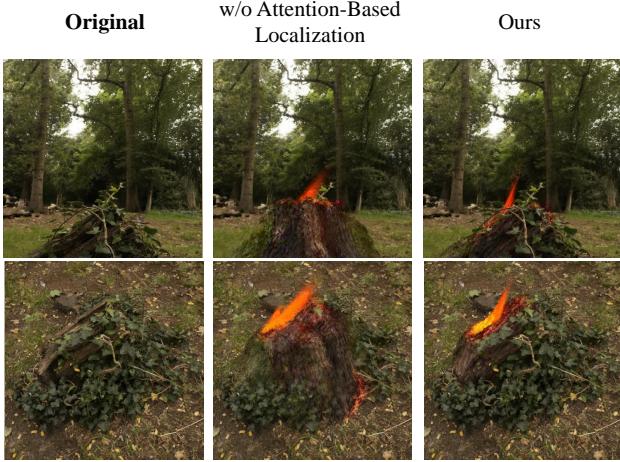
**Figure 8:** *Ablation study on attention-based localization of editing areas in 3D. Prompt: "a * stump on fire".* *Utilizing attention-based localization in 3D allows for better preservation of the entire scene's details. Without our 3D localization, plants near the stump lose their form, whereas our complete method maintains the scene's similarity to the original scene from various perspectives.*
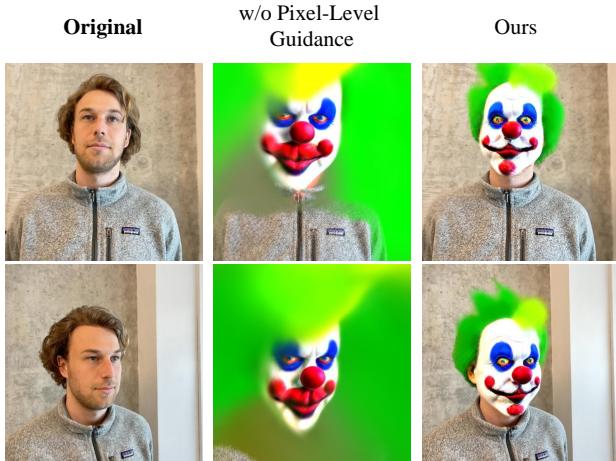


**Figure 9:** *Ablation study on Pixel-Level Guidance using pseudo-GT images. Prompt: "a * clown".* *Due to the strong fluidity of Gaussians, the background remains cluttered with 3D static localization when pixel-level guidance is removed. However, with the guidance of pseudo-GT images, our method effectively controls the Gaussians, enabling detailed editing.*

still messed up by the color from the hair of clown due the strong fluidity of Gaussians. However, when introducing pseudo-GT images for refinement, it can be visually observed from Figure 9 that it effectively reduces artifacts of the background and better maintains the consistency of non-editing areas before and after editing, thereby greatly improving the editing quality.

## 5. Conclusion

We propose GSEditPro, a novel text-based 3D scene editing framework capable of performing various editing operations. By leveraging the explicit nature of 3D Gaussian distributions, we have devised a method to add semantic labels based on attention maps to each Gaussian during the differentiable rendering process, achieving precise localization of editing regions. Additionally, we have fine-tuned the diffusion model to optimize Gaussians and devised an attention-based method to preserve details through pixel-level guidance. Our extensive experiments have demonstrated that the attention-based progressive localization in both 2D and 3D significantly enhances our framework, outperforming previous methods even with less prior information.

**Limitation.** Our method heavily relies on the generation ability of the 2D diffusion models, which may cause our 3D editing to fail when the 2D generation from the diffusion models is poor. As shown in Figure 10, with the 2D supervision generated by the diffusion model (in the middle column), it will be difficult to make satisfactory editing results that match the prompt (in the right column).



**Figure 10:** *A failed case. Prompt: "a plush * toy wearing shoes".* *The middle column shows a bad generation of diffusion models, leading to the failure of editing(the right column).*

## References

[BHE23] BROOKS T., HOLYNSKI A., EFROS A. A.: Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 18392–18402. 3, 6

[BMV*22] BARRON J. T., MILDENHALL B., VERBIN D., SRINIVASAN P. P., HEDMAN P.: Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 5470–5479. 6

[CBM*24] CHEN J.-K., BULÒ S. R., MÜLLER N., PORZI L., KONTSCHIEDER P., WANG Y.-X.: Consistdreamer: 3d-consistent 2d diffusion for high-fidelity scene editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2024), pp. 21071–21080. 3

[CCZ*23] CHEN Y., CHEN Z., ZHANG C., WANG F., YANG X., WANG Y., CAI Z., YANG L., LIU H., LIN G.: Gaussianeditor: Swift and controllable 3d editing with gaussian splatting. *arXiv preprint arXiv:2311.14521* (2023). 2, 3, 5, 6, 7, 8, 9

[CLV24] CHEN M., LAINA I., VEDALDI A.: Dge: Direct gaussian 3d editing by consistent multi-view editing. 3

[CZF*23] CEN J., ZHOU Z., FANG J., SHEN W., XIE L., JIANG D., ZHANG X., TIAN Q., ET AL.: Segment anything in 3d with nerfs. *Advances in Neural Information Processing Systems 36* (2023), 25971–25990. 3

[DYL*24] DI D., YANG J., LUO C., XUE Z., CHEN W., YANG X., GAO Y.: Hyper-3dg: Text-to-3d gaussian generation via hypergraph. *arXiv preprint arXiv:2403.09236* (2024). 3

[EKSX96] ESTER M., KRIEGEL H.-P., SANDER J., XU X.: Density-based spatial clustering of applications with noise. In *Int. Conf. knowledge discovery and data mining* (1996), vol. 240. 5

[FWZ*23] FANG J., WANG J., ZHANG X., XIE L., TIAN Q.: Gaussianeditor: Editing 3d gaussians delicately with text instructions. *arXiv preprint arXiv:2311.16037* (2023). 2, 3

[HMT*22] HERTZ A., MOKADY R., TENENBAUM J., ABERMAN K., PRITCH Y., COHEN-OR D.: Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626* (2022). 2, 5

[HSC*22] HO J., SAHARIA C., CHAN W., FLEET D. J., NOROUZI M., SALIMANS T.: Cascaded diffusion models for high fidelity image generation. *Journal of Machine Learning Research 23*, 47 (2022), 1–33. 2

[HTE*23] HAQUE A., TANCIK M., EFROS A. A., HOLYNSKI A., KANAZAWA A.: Instruct-nerf2nerf: Editing 3d scenes with instructions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 19740–19750. 2, 3, 6, 7, 8, 9

[JMB*22] JAIN A., MILDENHALL B., BARRON J. T., ABBEEL P., POOLE B.: Zero-shot text-guided object generation with dream fields. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2022), pp. 867–876. 3

[KKLD23] KERBL B., KOPANAS G., LEIMKÜHLER T., DRETTAKIS G.: 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics 42*, 4 (2023), 1–14. 2, 3, 6

[KMR*23] KIRILLOV A., MINTUN E., RAVI N., MAO H., ROLLAND C., GUSTAFSON L., XIAO T., WHITEHEAD S., BERG A. C., LO W.-Y., ET AL.: Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 4015–4026. 5

[KMS22] KOBAYASHI S., MATSUMOTO E., SITZMANN V.: Decomposing nerf for editing via feature field distillation. *Advances in Neural Information Processing Systems 35* (2022), 23311–23330. 4

[LGT*22] LIN C.-H., GAO J., TANG L., TAKIKAWA T., ZENG X., HUANG X., KREIS K., FIDLER S., LIU M.-Y., LIN T.-Y.: Magic3d: High-resolution text-to-3d content creation. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2022), 300–309. 2, 3

[LLF*23] LI Y., LIN Z.-H., FORSYTH D., HUANG J.-B., WANG S.: Climatenerf: Extreme weather synthesis in neural radiance field. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 3227–3238. 4

[LLSH23] LI J., LI D., SAVARESE S., HOI S.: BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML* (2023). 6

[LSZ*24] LI H., SHI H., ZHANG W., WU W., LIAO Y., WANG L., LEE L.-H., ZHOU P.: Dreamscene: 3d gaussian-based text-to-3d scene generation via formation pattern sampling. *arXiv preprint arXiv:2404.03575* (2024). 3

[LZZ*21] LIU S., ZHANG X., ZHANG Z., ZHANG R., ZHU J.-Y., RUSSELL B.: Editing conditional radiance fields. In *Proceedings of the IEEE/CVF international conference on computer vision* (2021), pp. 5773–5783. 3

[MPS*23] MIKAEILI A., PEREL O., SAFAEE M., COHEN-OR D., MAHDAVI-AMIRI A.: Sked: Sketch-guided text-based 3d editing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 14607–14619. 2, 3

[MRP*23] METZER G., RICHARDSON E., PATASHNIK O., GIRYES R., COHEN-OR D.: Latent-nerf for shape-guided generation of 3d shapes and textures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 12663–12673. 2, 3

[MST*21] MILDENHALL B., SRINIVASAN P. P., TANCIK M., BARRON J. T., RAMAMOORTHI R., NG R.: Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM 65*, 1 (2021), 99–106. 2, 3, 6

[PJBM22] POOLE B., JAIN A., BARRON J. T., MILDENHALL B.: Dreamfusion: Text-to-3d using 2d diffusion. *ArXiv abs/2209.14988* (2022). 2, 3, 5

[RBL*22] ROMBACH R., BLATTMANN A., LORENZ D., ESSER P., OMMER B.: High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2022), pp. 10684–10695. 2, 5

[RDN*22] RAMESH A., DHARIWAL P., NICHOL A., CHU C., CHEN M.: Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125 1*, 2 (2022), 3. 2

[RKH*21] RADFORD A., KIM J. W., HALLACY C., RAMESH A., GOH G., AGARWAL S., SASTRY G., ASKELL A., MISHKIN P., CLARK J., ET AL.: Learning transferable visual models from natural language supervision. In *International conference on machine learning* (2021), PMLR, pp. 8748–8763. 3, 6

[RKP*23] RAJ A., KAZA S., POOLE B., NIEMEYER M., RUIZ N., MILDENHALL B., ZADA S., ABERMAN K., RUBINSTEIN M., BARRON J., ET AL.: Dreambooth3d: Subject-driven text-to-3d generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 2349–2359. 3

[RLJ*23] RUIZ N., LI Y., JAMPANI V., PRITCH Y., RUBINSTEIN M., ABERMAN K.: 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* (2023), pp. 22500–22510. 2, 3, 5, 6, 8

[SCS*22] SAHARIA C., CHAN W., SAXENA S., LI L., WHANG J., DENTON E. L., GHASEMIPOUR K., GONTIJO LOPES R., KARAGOL AYAN B., SALIMANS T., ET AL.: Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems 35* (2022), 36479–36494. 2

[SF16] SCHONBERGER J. L., FRAHM J.-M.: Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2016), pp. 4104–4113. 5, 6

[SFHAE23] SELLA E., FIEBELMAN G., HEDMAN P., AVERBUCH-ELOR H.: Vox-e: Text-guided voxel editing of 3d objects. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 430–440. 3

[WCH*22] WANG C., CHAI M., HE M., CHEN D., LIAO J.: Clip-nerf: Text-and-image driven manipulation of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 3835–3844. 2, 3

[WDR*24] WANG X., DARRELL T., RAMBHATLA S. S., GIRDHAR R., MISRA I.: Instancediffusion: Instance-level control for image generation. *arXiv preprint arXiv:2402.03290* (2024). 2

[WLW*24] WANG Z., LU C., WANG Y., BAO F., LI C., SU H., ZHU J.: Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *Advances in Neural Information Processing Systems 36* (2024). 2

[WYW*24] WANG Y., YI X., WU Z., ZHAO N., CHEN L., ZHANG H.: View-consistent 3d editing with gaussian splatting. *arXiv preprint arXiv:2403.11868* (2024). 3

[XCC*24] XU T., CHEN J., CHEN P., ZHANG Y., YU J., YANG W.: Tiger: Text-instructed 3d gaussian retrieval and coherent editing. *arXiv preprint arXiv:2405.14455* (2024). 3

[XH22] XU T., HARADA T.: Deforming radiance fields with cages. In *European Conference on Computer Vision* (2022), Springer, pp. 159–175. 4

[YDYK23] YE M., DANELLJAN M., YU F., KE L.: Gaussian grouping: Segment and edit anything in 3d scenes. *arXiv preprint arXiv:2312.00732* (2023). 2, 3

[YSL*22] YUAN Y.-J., SUN Y.-T., LAI Y.-K., MA Y., JIA R., GAO L.: Nerf-editing: geometry editing of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 18353–18364. 4

[YXK*22] YU J., XU Y., KOH J. Y., LUONG T., BAID G., WANG Z., VASUDEVAN V., KU A., YANG Y., AYAN B. K., ET AL.: Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789 2*, 3 (2022), 5. 2

[ZKC*24] ZHUANG J., KANG D., CAO Y.-P., LI G., LIN L., SHAN Y.: Tip-editor: An accurate 3d editor following both text-prompts and image-prompts. *arXiv preprint arXiv:2401.14828* (2024). 2

[ZRA23] ZHANG L., RAO A., AGRAWALA M.: Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 3836–3847. 2

[ZRX*24] ZHOU X., RAN X., XIONG Y., HE J., LIN Z., WANG Y., SUN D., YANG M.-H.: Gala3d: Towards text-to-3d complex scene generation via layout-guided generative gaussian splatting. *arXiv preprint arXiv:2402.07207* (2024). 3

[ZWL*23] ZHUANG J., WANG C., LIN L., LIU L., LI G.: Dreameditor: Text-driven 3d scene editing with neural fields. In *SIGGRAPH Asia 2023 Conference Papers* (2023), pp. 1–10. 3, 6, 8, 9