

Feature 3DGS: Supercharging 3D Gaussian Splatting to Enable Distilled Feature Fields

Shijie Zhou¹ Haoran Chang^{1*} Sicheng Jiang^{1*} Zhiwen Fan² Zehao Zhu² Dejia Xu²
Pradyumna Chari¹ Suya You³ Zhangyang Wang² Achuta Kadambi¹

¹University of California, Los Angeles ²University of Texas at Austin ³DEVCOM Army Research Laboratory

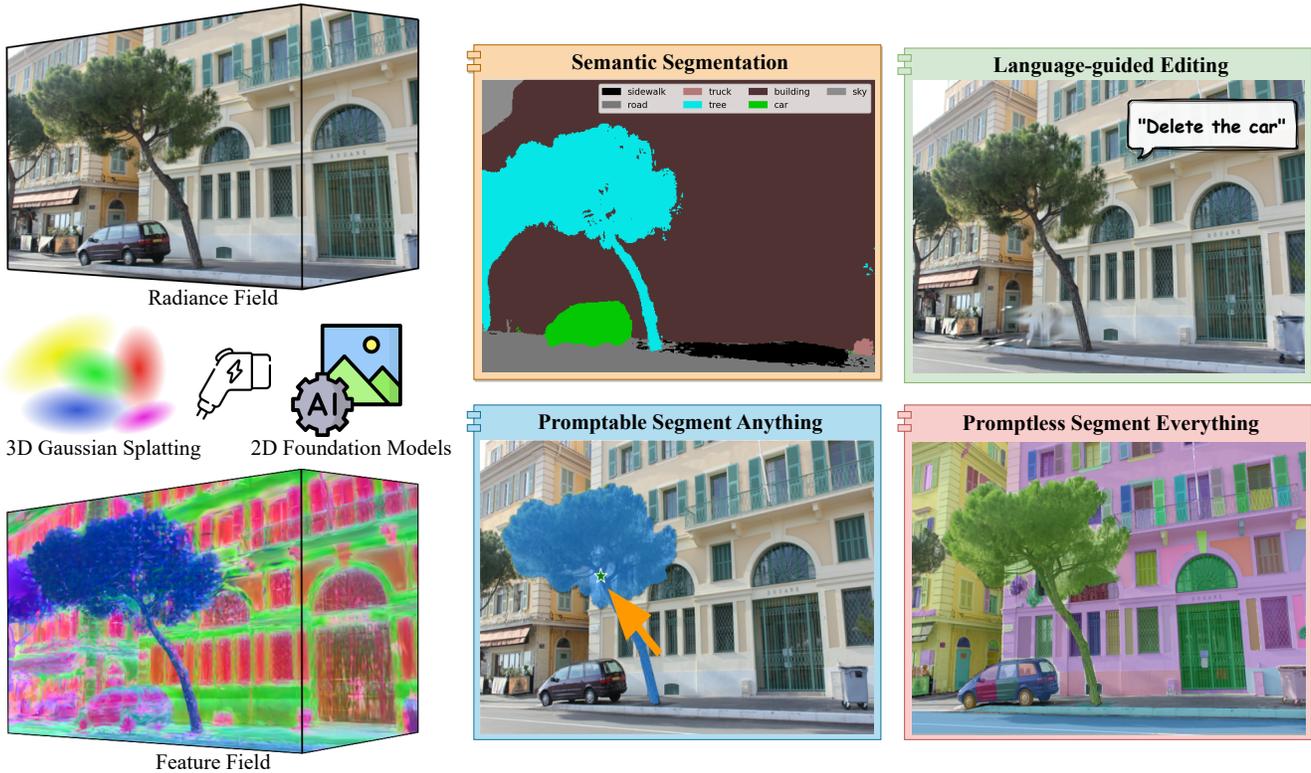


Figure 1. **Feature 3DGS**. We present a general method that significantly enhances 3D Gaussian Splatting through the integration of large 2D foundation models via feature field distillation. This advancement extends the capabilities of 3D Gaussian Splatting beyond mere novel view synthesis. It now encompasses a range of functionalities, including semantic segmentation, language-guided editing, and promptable segmentations such as “segment anything” or automatic segmentation of everything from any novel view. Scene from [16].

Abstract

3D scene representations have gained immense popularity in recent years. Methods that use Neural Radiance fields are versatile for traditional tasks such as novel view synthesis. In recent times, some work has emerged that aims to extend the functionality of NeRF beyond view synthesis, for semantically aware tasks such as editing and segmentation using 3D feature field distillation from 2D foundation

models. However, these methods have two major limitations: (a) they are limited by the rendering speed of NeRF pipelines, and (b) implicitly represented feature fields suffer from continuity artifacts reducing feature quality. Recently, 3D Gaussian Splatting has shown state-of-the-art performance on real-time radiance field rendering. In this work, we go one step further: in addition to radiance field rendering, we enable 3D Gaussian splatting on arbitrary-dimension semantic features via 2D foundation model distillation. This translation is not straightforward: naively in-

*Equal contribution.

corporating feature fields in the 3DGS framework encounters significant challenges, notably the disparities in spatial resolution and channel consistency between RGB images and feature maps. We propose architectural and training changes to efficiently avert this problem. Our proposed method is general, and our experiments showcase novel view semantic segmentation, language-guided editing and segment anything through learning feature fields from state-of-the-art 2D foundation models such as SAM and CLIP-LSeg. Across experiments, our distillation method is able to provide comparable or better results, while being significantly faster to both train and render. Additionally, to the best of our knowledge, we are the first method to enable point and bounding-box prompting for radiance field manipulation, by leveraging the SAM model. Project website at: <https://feature-3dgs.github.io/>.

1. Introduction

3D scene representation techniques have been at the forefront of computer vision and graphics advances in recent years. Methods such as Neural Radiance Fields (NeRFs) [30], and works that have followed up on it, have enabled learning implicitly represented 3D fields that are supervised on 2D images using the rendering equation. These methods have shown great promise for tasks such as novel view synthesis. However, since the implicit function is only designed to store local radiance information at every 3D location, the information contained in the field is limited from the perspective of downstream applications.

More recently, NeRF-based methods have attempted to use the 3D field to store additional descriptive features for the scene, in addition to the radiance [10, 19, 21, 46]. These features, when rendered into feature images, can then provide additional semantic information for the scene, enabling downstream tasks such as editing, segmentation and so on. However, feature field distillation through such a method is subject to a major disadvantage: NeRF-based methods can be natively slow to train as well as to infer. This is further complicated by model capacity issues: if the implicit representation network is kept fixed, while requiring it to learn an additional feature field (to not make the rendering and inference speeds even slower), the quality of the radiance field, as well as the feature field is likely to be affected unless the weight hyperparameter is meticulously tuned [21].

A recent alternative for implicit radiance field representations is the 3D Gaussian splatting-based radiance field proposed by Kerbl et al. [18]. This explicitly-represented field using 3D Gaussians is found to have superior training speeds and rendering speeds when compared with NeRF-based methods, while retaining comparable or better quality of rendered images. This speed of rendering while retaining high quality has paved the way for real-time ren-

dering applications, such as in VR and AR, that were previously found to be difficult. However, the 3D Gaussian splatting framework suffers the same representation limitation as NeRFs: natively, the framework does not support joint learning of semantic features and radiance field information at each Gaussian.

In this work, we present Feature 3DGS: the first feature field distillation technique based on the 3D Gaussian Splatting framework. Specifically, we propose learning a semantic feature at each 3D Gaussian, in addition to color information. Then, by splatting and rasterizing the feature vectors differentiably, the distillation of the feature field is possible using guidance from 2D foundation models. While the structure is natural and simple, enabling fast yet high-quality feature field distillation is not trivial: as the dimension of the learnt feature at each Gaussian increases, both training and rendering speeds drop drastically. We therefore propose learning a structured lower-dimensional feature field, which is later upsampled using a lightweight convolutional decoder at the end of the rasterization process. Therefore, this pipeline enables us to achieve improved feature field distillation at faster training and rendering speeds than NeRF-based methods, enabling a range of applications, including semantic segmentation, language-guided editing, promptable/promptless instance segmentation and so on.

In summary, our contributions are as follows:

- A novel 3D Gaussian splatting inspired framework for feature field distillation using guidance from 2D foundation models.
- A general distillation framework capable of working with a variety of feature fields such as CLIP-LSeg, Segment Anything (SAM) and so on.
- Up to $2.7\times$ faster feature field distillation and feature rendering over NeRF-based method by leveraging low-dimensional distillation followed by learnt convolutional upsampling.
- Up to **23%** improvement on mIoU for tasks such as semantic segmentation.

2. Related Work

2.1. Implicit Radiance Field Representations

Implicit neural representations have achieved remarkable success in recent years across a variety of areas within the field of computer vision [1, 28, 30, 32, 34, 48]. NeRF [30] demonstrates outstanding performance in novel view synthesis by representing 3D scenes with a coordinate-based neural network. In mip-NeRF [1], point-based ray tracing is replaced using cone tracing to combat aliasing. Zip-NeRF [2] utilized an anti-aliased grid-based technique to boost the radiance field performance. Instant-NGP [31] reduces the cost for neural primitives with a versatile new

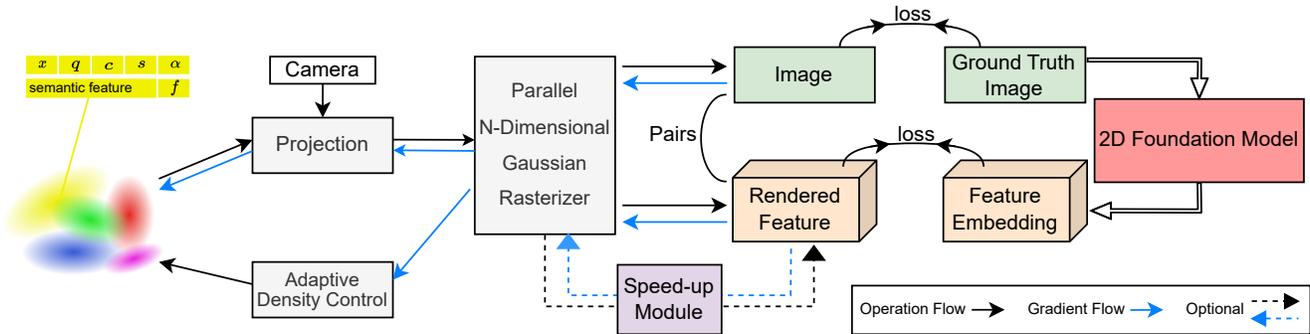


Figure 2. **An overview of our method.** We adopt the same 3D Gaussian initialization from sparse SfM point clouds as utilized in 3DGS, with the addition of an essential attribute: the *semantic feature*. Our primary innovation lies in the development of a Parallel N-dimensional Gaussian Rasterizer, complemented by a convolutional speed-up module as an optional branch. This configuration is adept at rapidly rendering arbitrarily high-dimensional features without sacrificing downstream performance.

input encoding that permits the use of a smaller network without sacrificing quality, thus significantly reducing the number of floating point and memory access operations. IBRNet [47], MVSNeRF [6], and PixelNeRF [52] construct a generalizable 3D representation by leveraging features gathered from various observed viewpoints. However, NeRF-based methods are hindered by slow rendering speeds and substantial memory usage during training, a consequence of their implicit design.

2.2. Explicit Radiance Field Representations

Pure implicit radiance fields are slow to operate and usually require millions of times querying a neural network for rendering a large-scale scene. Marrying explicit representations into implicit radiance fields enjoys the best of both worlds. Triplane [5], TensorRF [7], K-Plane [12], TILED [51] adopt tensor factorization to obtain efficient explicit representation. InstantNGP [31] utilizes multi-scale hash grids to work with large-scale scenes. BlockNeRF [44] further extends NeRF to render city-scale scenes spanning multiple blocks. Point NeRF [49] uses neural 3D points for representing and rendering a continuous radiance volume. NU-MCC [25] similarly utilizes latent point features but focuses on shape completion tasks. Unlike NeRF-style volumetric rendering, 3D Gaussian Splatting introduces point-based α -blending and an efficient point-based rasterizer. Our work follows 3D Gaussians Splatting, where we represent the scene using explicit point-based 3D representation, *i.e.* anisotropic 3D Gaussians.

2.3. Feature Field Distillation

Enabling simultaneously novel view synthesis and representing feature fields is well explored under NeRF [30] literature. Pioneering works such as Semantic NeRF [53] and Panoptic Lifting [41] have successfully embedded se-

semantic data from segmentation networks into 3D spaces. Their research has shown that merging noisy or inconsistent 2D labels in a 3D environment can yield sharp and precise 3D segmentation. Further extending this idea, techniques like those presented in [37] have demonstrated the effectiveness of segmenting objects in 3D with minimal user input, like rudimentary foreground-background masks. Beyond optimizing NeRF with estimated labels, Distilled Feature Fields [21], NeRF-SOS [10], LERF [19], and Neural Feature Fusion Fields [46] have delved into embedding pixel-aligned feature vectors from technologies such as LSeg or DINO [4] into NeRF frameworks. Additionally, [13, 26, 27, 39, 40, 45, 50] also explore feature fusion and manipulation in 3D. Feature 3DGS shares a similar idea for distilling 2D well-trained models, but also demonstrates an effective way of distilling into explicit point-based 3D representations, for simultaneous photo-realistic view synthesis and label map rendering.

3. Method

NeRF-based feature field distillation, as explored in [21], utilizes two distinct branches of MLPs to output the color c and feature f . Subsequently, the RGB image and high-dimensional feature map are rendered individually through volumetric rendering. The transition from NeRF to 3DGS is not as straightforward as simply rasterizing RGB images and feature maps independently. Typically, feature maps have fixed dimensions that often differ from those of RGB images. Due to the tile-based rasterization procedure and shared attributes between images and feature maps, rendering them independently can be problematic. A naive approach is to adopt a two-stage training method that rasterizes them separately. However, this approach could result in suboptimal quality for both RGB images and feature maps, given the high-dimensional correlations of semantic

features with the shared attributes of RGB.

In this section, we introduce a novel pipeline for high-dimensional feature rendering and feature field distillation, which enables 3D Gaussians to explicitly represent both radiance fields and feature fields. Our proposed parallel N -dimensional Gaussian rasterizer and speed-up module can effectively solve the aforementioned problems and is capable of rendering arbitrary dimensional semantic feature map. An overview of our method is shown in Fig. 2. Our proposed method is general and compatible with any 2D foundation model, by distilling the semantic features into a 3D feature field using 3D Gaussian splatting. In our experiments, we employ SAM [20] and LSeg [23], facilitating promptable, promptless (zero-shot [3] [14]) and language-driven computer vision tasks in a 3D context.

3.1. High-dimensional Semantic Feature Rendering

To develop a general feature field distillation pipeline, our method should be able to render 2D feature maps of arbitrary size and feature dimension, in order to cope with different kinds of 2D foundation models. To achieve this, we use the rendering pipeline based on the differentiable Gaussian splatting framework proposed by [18] as our foundation. We follow the same 3D Gaussians initialization technique using Structure from Motion [38]. Given this initial point cloud, each point $x \in \mathbb{R}^3$ within it can be described as the center of a Gaussian. In world coordinates, the 3D Gaussians are defined by a full 3D covariance matrix Σ , which is transformed to Σ' in camera coordinates when 3D Gaussians are projected to 2D image / feature map space [55]:

$$\Sigma' = JW\Sigma W^T J^T, \quad (1)$$

where W is the world-to-camera transformation matrix and J is the Jacobian of the affine approximation of the projective transformation. Σ is physically meaningful only when it is positive semi-definite — a condition that cannot always be guaranteed during optimization. This issue can be addressed by decomposing Σ into rotation matrix R and scaling matrix S :

$$\Sigma = RSS^T R^T, \quad (2)$$

Practically, the rotation matrix R and the scaling matrix S are stored as a rotation quaternion $q \in \mathbb{R}^4$ and a scaling factor $s \in \mathbb{R}^3$ respectively. Besides the aforementioned optimizable parameters, an opacity value $\alpha \in \mathbb{R}$ and spherical harmonics (SH) up to the 3rd order are also stored in the 3D Gaussians. In practice, we optimize the zeroth-order SH for the first 1000 iterations, which equates to a simple diffuse color representation $c \in \mathbb{R}^3$, and we introduce 1 band every 1000 iterations until all 4 bands of SH are represented. Additionally, we incorporate the semantic feature $f \in \mathbb{R}^N$, where N can be any arbitrary number representing the latent dimension of the feature. In summary, for the

i -th 3D Gaussian, the optimizable attributes are given by $\Theta_i = \{x_i, q_i, s_i, \alpha_i, c_i, f_i\}$.

Upon projecting the 3D Gaussians into a 2D space, the color C of a pixel and the feature value F_s of a feature map pixel are computed by volumetric rendering which is performed using front-to-back depth order [22]:

$$C = \sum_{i \in \mathcal{N}} c_i \alpha_i T_i, \quad F_s = \sum_{i \in \mathcal{N}} f_i \alpha_i T_i, \quad (3)$$

where $T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$, \mathcal{N} is the set of sorted Gaussians overlapping with the given pixel, T_i is the transmittance, defined as the product of opacity values of previous Gaussians overlapping the same pixel. The subscript s in F_s denotes ‘‘student’’, indicating that this rendered feature is per-pixel supervised by the ‘‘teacher’’ feature F_t . The latter represents the latent embedding obtained by encoding the ground truth image using the encoder of 2D foundation models. This supervisory relationship underscores the instructional dynamic between F_s and F_t in our model. In essence, our approach involves distilling [17] the large 2D teacher model into our small 3D student explicit scene representation model through differentiable volumetric rendering.

In the rasterization stage, we adopted a joint optimization method, as opposed to rasterizing the RGB image and feature map independently. Both image and feature map utilize the same tile-based rasterization procedure, where the screen is divided into 16×16 tiles, and each thread processes one pixel. Subsequently, 3D Gaussians are culled against both the view frustum and each tile. Owing to their shared attributes, both the feature map and RGB image are rasterized to the same resolution but in different dimensions, corresponding to the dimensions of c_i and f_i initialized in the 3D Gaussians. This approach ensures that the fidelity of the feature map is rendered as high as that of the RGB image, thereby preserving per-pixel accuracy.

3.2. Optimization and Speed-up

The loss function is the photometric loss combined with the feature loss:

$$\mathcal{L} = \mathcal{L}_{rgb} + \gamma \mathcal{L}_f, \quad (4)$$

with

$$\begin{aligned} \mathcal{L}_{rgb} &= (1 - \lambda) \mathcal{L}_1(I, \hat{I}) + \lambda \mathcal{L}_{D-SSIM}(I, \hat{I}), \\ \mathcal{L}_f &= \|F_t(I) - F_s(\hat{I})\|_1. \end{aligned}$$

where I is the ground truth image and \hat{I} is our rendered image. The latent embedding $F_t(I)$ is derived from the 2D foundation model by encoding the image I , while $F_s(\hat{I})$ represents our rendered feature map. To ensure identical resolution $H \times W$ for the per-pixel \mathcal{L}_1 loss calculation, we apply bilinear interpolation to resize $F_s(\hat{I})$ accordingly. In

practice, we set the weight hyperparameters $\gamma = 1.0$ and $\lambda = 0.2$.

It is important to note that in NeRF-based feature field distillation, the scene is implicitly represented as a neural network. In this configuration, as discussed in [21], the branch dedicated to the feature field shares some layers with the radiance field. This overlap could potentially lead to interference, where learning the feature fields might adversely affect the radiance fields. To address this issue, a compromise approach is to set γ to a low value, meaning the weight of the feature field is much smaller than that of the radiance field during the optimization. [21] also mentions that NeRF is highly sensitive to γ . Conversely, our explicit scene representation avoids this issue. Our equal-weighted joint optimization approach has demonstrated that the resulting high-dimensional semantic features significantly contribute to scene understanding and enhance the depiction of physical scene attributes, such as opacity and relative positioning. See the comparison between Ours and Base 3DGS in Tab. 1.

To optimize the semantic feature $f \in \mathbb{R}^N$, we minimize the difference between the rendered feature map $F_s(\hat{I}) \in \mathbb{R}^{H \times W \times N}$ and the teacher feature map $F_t(I) \in \mathbb{R}^{H \times W \times M}$, ideally with $N = M$. However, in practice, M tends to be a very large number due to the high latent dimensions in 2D foundation models (e.g. $M = 512$ for LSeg and $M = 256$ for SAM), making direct rendering of such high-dimensional feature maps time-consuming. To address this issue, we introduce a speed-up module at the end of the rasterization process. This module consists of a lightweight convolutional decoder that upsamples the feature channels with kernel size 1×1 . Consequently, it is feasible to initialize $f \in \mathbb{R}^N$ on 3D Gaussians with any arbitrary $N \ll M$ and to use this learnable decoder to match the feature channels. This allows us to not only effectively achieve $F_s(\hat{I}) \in \mathbb{R}^{H \times W \times M}$, but also significantly speed up the optimization process without compromising the performance on downstream tasks.

The advantages of implementing this convolutional speed-up module are threefold: Firstly, the input to the convolution layer, with a kernel size of 1×1 , is the resized rendered feature map, which is significantly smaller in size compared to the original image. This makes the 1×1 convolution operation computationally efficient. Secondly, this convolution layer is a learnable component, facilitating channel-wise communication within the high-dimensional rendered feature, enhancing the feature representation. Lastly, the module’s design is optional. Whether included or not, it does not impact the performance of downstream tasks, thereby maintaining the flexibility and adaptability of the entire pipeline.

3.3. Promptable Explicit Scene Representation

Foundation models provide a base layer of knowledge and skills that can be adapted for a variety of specific tasks and applications. We wish to use our feature field distillation approach to enable practical 3D representations of these features. Specifically, we consider two foundation models, namely Segment Anything [20], and LSeg [23]. The *Segment Anything Model (SAM)* [20] allows for both promptable and promptless zero-shot segmentation in 2D, without the need for specific task training. LSeg [23] introduces a language-driven approach to zero-shot semantic segmentation. Utilizing the image feature encoder with the DPT architecture [36] and text encoders from CLIP [35], LSeg extends text-image associations to a 2D pixel-level granularity. Through the teacher-student distillation, **our distilled feature fields facilitate the extension of all 2D functionalities — prompted by point, box, or text — into the 3D realm.**

Our promptable explicit scene representation works as follows: for a 3D Gaussian x among the N ordered Gaussians overlapping the target pixel, i.e. $x_i \in \mathcal{X}$ where $\mathcal{X} = \{x_1, \dots, x_N\}$, the activation score of a prompt τ on the 3D Gaussian x is calculated by cosine similarity between the query $q(\tau)$ in the feature space and the semantic feature $f(x)$ of the 3D Gaussian followed by a softmax:

$$s = \frac{f(x) \cdot q(\tau)}{\|f(x)\| \|q(\tau)\|}, \quad (5)$$

If we have a set \mathcal{T} of possible labels, such as a text label set for semantic segmentation or a point set of all the possible pixels for point-prompt, the probability of a prompt τ of a 3D Gaussian can be obtained by softmax:

$$\mathbf{p}(\tau|x) = \text{softmax}(s) = \frac{\exp(s)}{\sum_{s_j \in \mathcal{T}} \exp(s_j)}. \quad (6)$$

We utilize the computed probabilities to filter out Gaussians with low probability scores. This selective approach enables various operations, such as extraction, deletion, or appearance modification, by updating the color $c(x)$ and opacity $\alpha(x)$ values as needed. With the newly updated color set $\{c_i\}_{i=1}^n$ and opacity set $\{\alpha_i\}_{i=1}^n$, where n is smaller than N , we can implement point-based α -blending to render the edited radiance field from any novel view.

4. Experiments

4.1. Novel view semantic segmentation

The number of classes of a dataset is usually limited from tens [9] to hundreds [54], which is insignificant to English words [24]. In light of the limitation, semantic features empower models to comprehend unseen labels by mapping

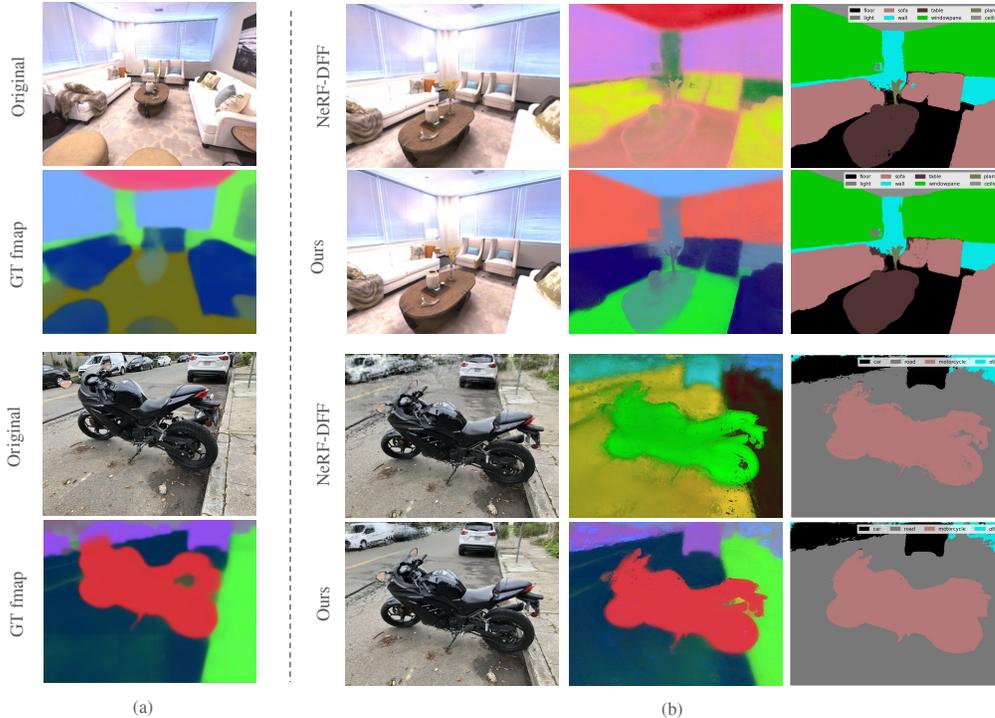


Figure 3. **Novel view semantic segmentation (Lseg) results on scenes from Replica dataset [43] and LLFF dataset [29].** (a) We show examples of original images in training views together with the ground-truth feature visualizations. (b) We compare the qualitative segmentation results using our Feature 3DGS with the NeRF-DFF [21]. Our inference is **1.66** \times faster when rendered feature $dim = 128$. Our method demonstrates more fine-grained segmentation results with higher-quality feature maps.

Metrics	PSNR(\pm s.d.) \uparrow	SSIM(\pm s.d.) \uparrow	LPIPS(\pm s.d.) \downarrow
Ours (w/ speed-up)	37.012 (\pm 0.07)	0.971 (\pm 5.3e-4)	0.023 (\pm 2.9e-4)
Ours	36.915 (\pm 0.05)	0.970 (\pm 5.7e-4)	0.024 (\pm 1.1e-3)
Base 3DGS	36.133 (\pm 0.06)	0.965 (\pm 1.5e-4)	0.033 (\pm 1.2e-3)

Table 1. **Performance on Replica Dataset.** (average performance for 5K training iterations, speed-up module rendered feature $dim = 128$). Boldface font represents the preferred results.

Metrics	mIoU \uparrow	accuracy \uparrow	FPS \uparrow
Ours (w/ speed-up)	0.782	0.943	14.55
Ours	0.787	0.943	6.84
NeRF-DFF	0.636	0.864	5.38

Table 2. **Performance of semantic segmentation on Replica dataset compared to NeRF-DFF.** (speed-up module rendered feature $dim = 128$). Boldface font represents the preferred results.

semantically close labels to similar regions in the embedding space, as articulated by Li et al [23]. This advancement notably promotes the scalability in information acquisition and scene understanding, facilitating a profound

comprehension of intricate scenes. We distill LSeg feature for this novel view semantic segmentation task. Our experiments demonstrate the improvement of incorporating semantic feature over the naive 3D Gaussian rasterization method [18]. In Tab. 1, we show that our model surpasses the baseline 3D Gaussian model in performance metrics on Replica dataset [43] with 5000 training iterations for all three models. Noticeably, the integration of the speed-up module to our model does not compromise the performance.

In our further comparison with NeRF-DFF [21] using the Replica dataset, we address the potential trade-off between the quality of the semantic feature map and RGB images. In Tab. 2, our model demonstrates higher accuracy and mean intersection-over-union (mIoU). Additionally, by incorporating our speed-up module, we achieved more than double the frame rate per second (FPS) of our full model while maintaining comparable performance. In Fig. 3 (b) the last column, our approach yields better visual quality on novel views and semantic segmentation masks for both synthetic and real scenes compared to NeRF-DFF.

4.2. Segment Anything from Any View

SAM excels in performing precise instance segmentation, utilizing interactive points and boxes as prompts to auto-

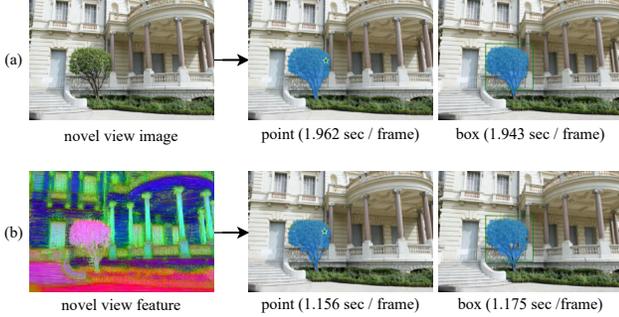


Figure 4. **Comparison of SAM segmentation results obtained by** (a) naively applying the SAM encoder-decoder module to a novel-view rendered image **with** (b) directly decoding a rendered feature. Our method is up to $1.7\times$ faster in total inference speed including rendering and segmentation while preserving the quality of segmentation masks. Scene from [16].

matically segment objects in any 2D image. In our experiments, we extend this capability to 3D, aiming to achieve fast and accurate segmentation from any viewpoint. Our distilled feature field enables the model to render the SAM feature map directly for any given camera pose. As such, the SAM decoder is the only component needed to interact with the input prompt and produce the segmentation mask, thereby bypassing the need to synthesize a novel view image first and then process it through the entire SAM encoder-decoder pipeline. Furthermore, to enhance training and inference speed, we use the speed-up module in this experiment. In practice, we set the rendered feature dimension to 128, which is half of SAM’s latent dimension of 256, maintaining the comparable quality of segmentation.

In Fig. 4, we compare the results of both point and box prompted segmentation on novel views using the naive approach (SAM encoder + decoder) and our proposed feature field approach (SAM decoder only). We achieve nearly equivalent segmentation quality, but our method is up to $1.7\times$ faster. In Fig. 5, we contrast our method with NeRF-DFF [21]. Our rendered features not only yield higher quality mask boundaries, as evidenced by the bear’s leg, but also deliver more accurate and comprehensive instance segmentation, capable of segmenting finer-grained instances (as illustrated by the more ‘detailed’ mask on the far right). Additionally, we use PCA-based feature visualization [33] to demonstrate that our high-quality segmentation masks result from superior feature rendering.

4.3. Language-guided Editing

In this section, we showcase the capability of our Feature 3DGS, distilled from LSeg, to perform editable novel view synthesis. The process begins by querying the feature field with a text prompt, typically comprising an edit operation followed by a target object, such as “extract the

car”. For text encoding, we employ a ViT-B/32 CLIP encoder. Our editing pipeline capitalizes on the semantic features queried from the 3D feature field, rather than relying on a 2D rendered feature map. We compute semantic scores for each 3D Gaussian, represented by a K -dimensional vector (where K is the number of object categories), using a softmax function. Subsequently, we engage in either soft selection (by setting a threshold value) or hard selection (by filtering based on the highest score across K categories). To identify the region for editing, we generate a “Gaussian mask” through thresholding on the score matrix, which is then utilized for modifications to color c and opacity α on 3D Gaussians.

In Fig. 6, we showcase our novel view editing results, achieved through various operations prompted by language inputs. Specifically, we conduct editing tasks such as extraction, deletion, and appearance modification on text-specified targets within diverse scenes. As illustrated in Fig. 6 (a), we successfully extract an entire banana from the scene. Notably, by leveraging 3D Gaussians to update the rendering parameters, our Feature 3DGS model gains an understanding of the 3D scene environment from any viewpoint. This enables the model to reconstruct occluded or invisible parts of the scene, as evidenced by the complete extraction of a banana initially hidden by an apple in view 1. Furthermore, compared with our edit results with NeRF-DFF, our method stands out by providing a cleaner extraction with little floaters in the background. Additionally, in Fig. 6 (b), our model is able to delete objects like cars while retaining the background elements, such as plants, due to the opacity updates in 3DGS, showcasing its 3D scene awareness. Moreover, we also demonstrate the model’s capability in modifying the appearance of specific objects, like ‘sidewalk’ and ‘leaves’, without affecting adjacent objects’ appearance (e.g., the ‘stop sign’ remains red).

5. Discussion and Conclusion

In this work, we present a notable advancement in explicit 3D scene representation by integrating 3D Gaussian Splatting with feature field distillation from 2D foundation models, a development that not only broadens the scope of radiance fields beyond traditional uses but also addresses key limitations of previous NeRF-based methods in implicitly represented feature fields. Our work, as showcased in various experiments including complex semantic tasks like editing, segmentation, and language-prompted interactions with models like CLIP-LSeg and SAM, opening the door to a brand new semantic, editable, and promptable explicit 3D scene representation.

However, our Feature 3DGS framework does have its inherent limitations. The student feature’s limited access to the ground truth feature restricts the overall performance, and the imperfections of the teacher network further con-

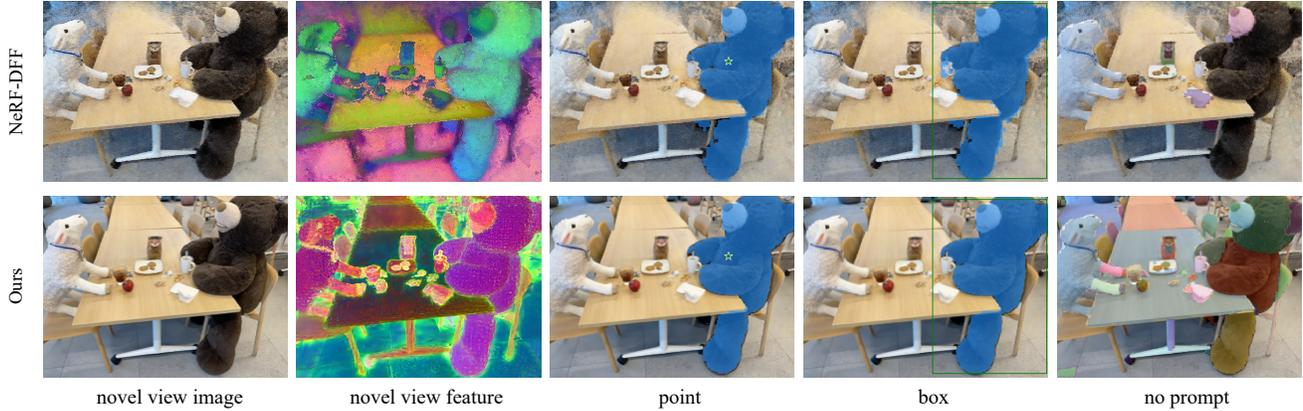


Figure 5. **Novel view segmentation (SAM) results compared with NeRF-DFF.** (Upper) NeRF-DFF method presents lower-quality segmentation masks - note the failure on segmenting the cup from the bear and the coarse-grained mask boundary on the bear’s leg in box-prompted results. (Lower) Our method provides higher-quality masks with more fine-grained segmentation details. Scene from [19].

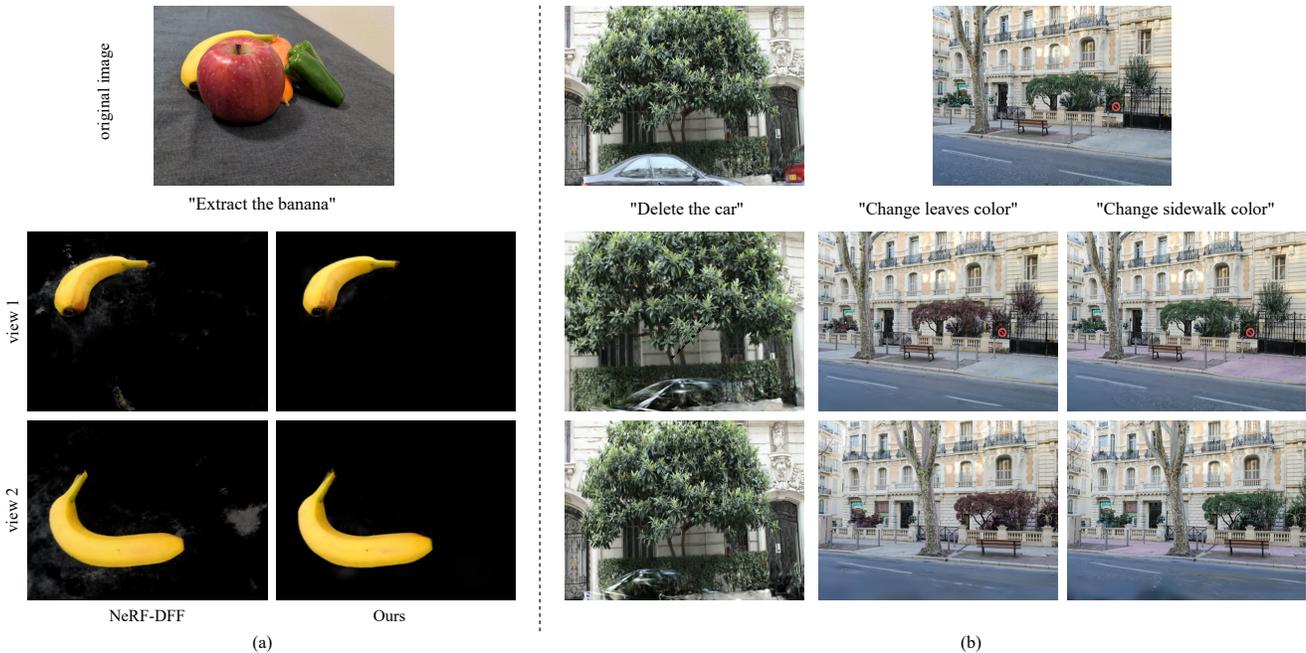


Figure 6. **Demonstration of results with various language-guided edit operations by querying the 3D feature field and comparison with NeRF-DFF** (a) We compare our edit results with NeRF-DFF method on the sample dataset provided by NeRF-DFF [21]. Note that our method outperforms NeRF-DFF method by extracting the entire banana hidden by an apple in the original image and with less floaters in the background. (b) We demonstrate results with deletion and appearance modification on different targets. Note that the car is deleted with background preserved, and the appearance of the leaves changes with the appearance of the stop sign remained the same.

strain our framework’s effectiveness. In addition, our adaptation of the original 3DGS pipeline, which inherently generates noise-inducing floaters, poses another challenge, affecting our model’s optimal performance.

Acknowledgement

We thank the Visual Machines Group (VMG) at UCLA and Visual Informatics Group at UT Austin (VITA) for feedback

and support. This project was supported by the US DoD LUCI (Laboratory University Collaboration Initiative) Fellowship and partially supported by ARL grants W911NF-20-2-0158 and W911NF-21-2-0104 under the cooperative A2I2 program. Z.W. is partially supported by the ARL grant W911NF2120064 under the cooperative A2I2 program, and an Army Young Investigator Award. A.K. is supported by a DARPA Young Faculty Award, NSF CAREER Award, and Army Young Investigator Award.

References

- [1] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. *ICCV*, 2021. 2
- [2] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Zip-nerf: Anti-aliased grid-based neural radiance fields. *ICCV*, 2023. 2
- [3] Maxime Bucher, Tuan-Hung Vu, Matthieu Cord, and Patrick Pérez. Zero-shot semantic segmentation. *Advances in Neural Information Processing Systems*, 32, 2019. 4
- [4] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 3
- [5] Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16123–16133, 2022. 3
- [6] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14124–14133, 2021. 3
- [7] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision*, pages 333–350. Springer, 2022. 3
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 13
- [9] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111:98–136, 2015. 5
- [10] Zhiwen Fan, Peihao Wang, Yifan Jiang, Xinyu Gong, De-jia Xu, and Zhangyang Wang. Nerf-sos: Any-view self-supervised object segmentation on complex scenes. *arXiv preprint arXiv:2209.08776*, 2022. 2, 3
- [11] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022. 12
- [12] Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12479–12488, 2023. 3
- [13] Rahul Goel, Dhawal Sirikonda, Saurabh Saini, and PJ Narayanan. Interactive segmentation of radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4201–4211, 2023. 3
- [14] Zhangxuan Gu, Siyuan Zhou, Li Niu, Zihan Zhao, and Liqing Zhang. Context-aware feature generation for zero-shot semantic segmentation. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 1921–1929, 2020. 4
- [15] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16000–16009, 2022. 13
- [16] Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel Brostow. Deep blending for free-viewpoint image-based rendering. *ACM Transactions on Graphics (ToG)*, 37(6):1–15, 2018. 1, 7
- [17] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 4
- [18] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (ToG)*, 42(4):1–14, 2023. 2, 4, 6, 12
- [19] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lrf: Language embedded radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19729–19739, 2023. 2, 3, 8
- [20] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023. 4, 5, 13
- [21] Sosuke Kobayashi, Eiichi Matsumoto, and Vincent Sitzmann. Decomposing nerf for editing via feature field distillation. In *NeurIPS*, 2022. 2, 3, 5, 6, 7, 8, 14
- [22] Georgios Kopanas, Julien Philip, Thomas Leimkühler, and George Drettakis. Point-based neural rendering with per-view optimization. In *Computer Graphics Forum*, pages 29–43. Wiley Online Library, 2021. 4
- [23] Boyi Li, Kilian Q. Weinberger, Serge Belongie, Vladlen Koltun, and René Ranftl. Language-driven semantic segmentation, 2022. 4, 5, 6
- [24] Xiang Li, Tianhan Wei, Yau Pun Chen, Yu-Wing Tai, and Chi-Keung Tang. Fss-1000: A 1000-class dataset for few-shot segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2869–2878, 2020. 5
- [25] Stefan Lionar, Xiangyu Xu, Min Lin, and Gim Hee Lee. Numcc: Multiview compressive coding with neighborhood decoder and repulsive udf. *arXiv preprint arXiv:2307.09112*, 2023. 3
- [26] Kunhao Liu, Fangneng Zhan, Jiahui Zhang, Muyu Xu, Yingchen Yu, Abdulmotaleb El Saddik, Christian Theobalt, Eric Xing, and Shijian Lu. Weakly supervised 3d open-vocabulary segmentation. *Advances in Neural Information Processing Systems*, 36, 2024. 3

- [27] Kirill Mazur, Edgar Sucar, and Andrew J Davison. Feature-realistic neural fusion for real-time, open set scene understanding. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 8201–8207. IEEE, 2023. 3
- [28] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4460–4470, 2019. 2
- [29] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 38(4):1–14, 2019. 6, 13, 17
- [30] 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. 2, 3
- [31] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (TOG)*, 41(4):1–15, 2022. 2, 3
- [32] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 165–174, 2019. 2
- [33] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830, 2011. 7, 13
- [34] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16*, pages 523–540. Springer, 2020. 2
- [35] 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. 5
- [36] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 12179–12188, 2021. 5
- [37] Zhongzheng Ren, Aseem Agarwala[†], Bryan Russell[†], Alexander G. Schwing[†], and Oliver Wang[†]. Neural volumetric object selection. In *CVPR*, 2022. ([†] alphabetic ordering). 3
- [38] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016. 4
- [39] Nur Muhammad Mahi Shafullah, Chris Paxton, Lerrel Pinto, Soumith Chintala, and Arthur Szlam. Clip-fields: Weakly supervised semantic fields for robotic memory. *arXiv preprint arXiv:2210.05663*, 2022. 3
- [40] William Shen, Ge Yang, Alan Yu, Jansen Wong, Leslie Pack Kaelbling, and Phillip Isola. Distilled feature fields enable few-shot language-guided manipulation. In *7th Annual Conference on Robot Learning*, 2023. 3
- [41] Yawar Siddiqui, Lorenzo Porzi, Samuel Rota Buló, Norman Müller, Matthias Nießner, Angela Dai, and Peter Kotschieder. Panoptic lifting for 3d scene understanding with neural fields. *arXiv preprint arXiv:2212.09802*, 2022. 3
- [42] Karl Stelzner, Kristian Kersting, and Adam R Kosiorok. De-composing 3d scenes into objects via unsupervised volume segmentation. *arXiv preprint arXiv:2104.01148*, 2021. 13, 14
- [43] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik Wijmans, Simon Green, Jakob J. Engel, Raul Mur-Artal, Carl Ren, Shobhit Verma, Anton Clarkson, Mingfei Yan, Brian Budge, Yajie Yan, Xiaqing Pan, June Yon, Yuyang Zou, Kimberly Leon, Nigel Carter, Jesus Briales, Tyler Gillingham, Elias Mueggler, Luis Pesqueira, Manolis Savva, Dhruv Batra, Hauke M. Strasdat, Renzo De Nardi, Michael Goesele, Steven Lovegrove, and Richard Newcombe. The replica dataset: A digital replica of indoor spaces, 2019. 6, 14
- [44] Matthew Tancik, Vincent Casser, Xinchun Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8248–8258, 2022. 3
- [45] Nikolaos Tsagkas, Oisín Mac Aodha, and Chris Xiaoxuan Lu. VI-fields: Towards language-grounded neural implicit spatial representations. In *2023 IEEE International Conference on Robotics and Automation*. IEEE, 2023. 3
- [46] Vadim Tschernezki, Iro Laina, Diane Larlus, and Andrea Vedaldi. Neural Feature Fusion Fields: 3D distillation of self-supervised 2D image representations. In *3DV*, 2022. 2, 3
- [47] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul Srinivasan, Howard Zhou, Jonathan T. Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *CVPR*, 2021. 3
- [48] Zhen Wang, Shijie Zhou, Jeong Joon Park, Despoina Paschalidou, Suya You, Gordon Wetzstein, Leonidas Guibas, and Achuta Kadambi. Alto: Alternating latent topologies for implicit 3d reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 259–270, 2023. 2
- [49] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. Point-nerf: Point-based neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5438–5448, 2022. 3

- [50] Jianglong Ye, Naiyan Wang, and Xiaolong Wang. Feature-nerf: Learning generalizable nerfs by distilling foundation models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8962–8973, 2023. 3
- [51] Brent Yi, Weijia Zeng, Sam Buchanan, and Yi Ma. Canonical factors for hybrid neural fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3414–3426, 2023. 3
- [52] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. 3
- [53] Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and Andrew Davison. In-place scene labelling and understanding with implicit scene representation. In *ICCV*, 2021. 3
- [54] Bolei Zhou, Hang Zhao, Xavier Puig, Tete Xiao, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Semantic understanding of scenes through the ade20k dataset. *International Journal of Computer Vision*, 127:302–321, 2019. 5
- [55] Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. Ewa volume splatting. In *Proceedings Visualization, 2001. VIS'01.*, pages 29–538. IEEE, 2001. 4

Feature 3DGS: Supercharging 3D Gaussian Splatting to Enable Distilled Feature Fields

Supplementary Material

This supplement is organized as follows:

- Section **A** contains network architecture details;
- Section **B** contains more details on the training and inference settings;
- Section **C** contains the details of the teacher features from 2D foundation models;
- Section **D** contains more details of Replica dataset experiment;
- Section **E** contains the algorithmic details of language-guided editing;
- Section **F** contains ablation studies of our method;
- Section **G** contains failure cases of complex scenes and reasoning analysis.

A. Details of Architectures

Parallel N-Dimensional Gaussian Rasterizer The parallel N-dimensional rasterizer maintains an architecture akin to the original 3DGS rasterizer. Moreover, it employs a point-based α -blending technique for rasterizing the feature map. To mitigate the issue of inconsistent spatial resolution inherent in tile-based rasterization, we ensure that both the RGB image and the feature map are rendered at matching sizes. Additionally, the parallel N-dimensional rasterizer is adaptable to various foundational models, implying that its dimensions are flexible and can vary accordingly. The detail of the Parallel N-Dimensional Gaussian Rasterization is in Algorithm 1.

Speed-up Module The primary objective of our Speed-up Module is to modify the feature map channels, enabling the relatively low-dimensional semantic features rendered from 3D Gaussians to align with the high-dimensional ground truth 2D feature map. To facilitate this, we employ a convolutional layer equipped with a 1×1 kernel, offering a direct and efficient solution. Given that we already possess the ground truth feature map from the teacher network, which serves as a target for the rendered feature map approximation, there is no necessity for a complex CNN architecture. This approach simplifies the process, ensuring effective feature alignment without the need for intricate feature extraction mechanisms. More experimental results regarding performance of the Speed-up Module are included in Sec. **F**.

B. Training and Inference Details

For the training and inference pipeline, one option is to directly render a feature map with the dimension same as the ground-truth feature (512 for LSeg encoding and 256 for

Algorithm 1 Parallel N-Dimensional Gaussian Rasterization

```
PointCloud  $\leftarrow$  Structure from Motion  $\triangleright$  Point Cloud
X, C  $\leftarrow$  PointCloud  $\triangleright$  Position, Colors
 $\Sigma, A, F$   $\leftarrow$  InitAttributes()
 $\triangleright$  Covariances, Opacities, Semantic Features
Ft(I)  $\leftarrow$  I applying Foundation Model  $\triangleright$  Feature Map
i  $\leftarrow$  0  $\triangleright$  Iteration Counter
repeat
  V, I, Ft  $\leftarrow$  GetTrainingView()
 $\triangleright$  Camera Pose, Image, Feature Map
   $\hat{I}, F_s$   $\leftarrow$  ParallelRasterizer(X, C,  $\Sigma, A, F, V$ )
 $\triangleright$  Rasterization
  L  $\leftarrow$  Loss(I,  $\hat{I}$ ) +  $\lambda$ Loss(Ft, Fs)
 $\triangleright$  Loss Calculation
  X,  $\Sigma, C, A, F$   $\leftarrow$  Adam(L)
 $\triangleright$  Backpropagation and Step
if IsRefinementStep(i) then
  for Gaussians(x, q, c,  $\alpha, f$ ) do
    if  $\alpha < \epsilon$  or IsTooLarge(x, q) then
      RemoveGaussian()
    end if
    if  $\nabla_p L > \tau_p$  then
      if  $\|S\| > \tau_S$  then
        SplitGaussian(x, q, c,  $\alpha, f$ )
 $\triangleright$  Over-reconstruction
      else
        CloneGaussian(x, q, c,  $\alpha, f$ )
 $\triangleright$  Under-reconstruction
      end if
    end if
  end for
end if
  i  $\leftarrow$  i + 1  $\triangleright$  Counter Increment
until Convergence
```

SAM encoding). Since rendering with such large dimension slows down the training, another option is to use our speed-up module: rendering a lower-dimensional feature map, which is later upsampled to the ground-truth feature dimension by a lightweight convolutional decoder. Similar to 3DGS [18], we use Adam optimizer for optimization during training and use a standard exponential decay scheduling similar to [11]. For image rendering, we mainly follow the 3DGS optimization strategy by using a 4 times lower image resolution and upsampling twice after 250 and

Dimension	8	16	32	64	128	256	512
Time	6:40	7:21	8:51	12:10	19:55	48:39	1:29:42
mIoU \uparrow	0.354	0.493	0.709	0.774	0.783	0.791	0.790
Accuracy \uparrow	0.735	0.880	0.927	0.939	0.944	0.944	0.943

Table A. **Evaluation of Semantic Segmentation Performance Across Different Dimensions.** This table presents the Time, mIoU, and Accuracy corresponding to each dimension level with LSeg feature.

Dimension	8	16	32	64	128	256	512
PSNR \uparrow	36.8879	36.8871	36.9671	36.9397	37.012	36.9474	36.9150
SSIM \uparrow	0.9699	0.9703	0.9706	0.9708	0.9706	0.9704	0.9703
LPIPS \downarrow	0.0234	0.0230	0.0226	0.0229	0.0228	0.0230	0.0236

Table B. **Evaluation of Image Quality Metrics Across Different Dimensions of Lseg feature.** This table presents the PSNR, SSIM, and LPIPS values corresponding to each dimension level with LSeg feature.

500 iterations. For feature rendering, we use Adam optimizer with a learning rate of $1e-3$. For the feature decoder network in the additional Speed-up Module, we use a separate Adam optimizer with a learning rate of $1e-4$.

C. Teacher Features

LSeg Feature For LSeg, we use CLIP ViT-L/16 image encoder for ground-truth feature preparation and ViT-L/16 text encoder for text encoding. The ground truth feature from the LSeg image encoder has feature size 360×480 with feature dimension 512. One can either choose to directly render a $h \times w$ feature with dimension 512 or use the Speed-up Module by rendering a lower-dimensional feature which is later upsampled back. In practice, we use rendered feature $dim = 128$ for Sec. 4.1 in our main paper.

To predict the semantic segmentation mask during inference, we reshape the rendered feature with shape (512, 360, 480) to $(360 \times 480, 512)$, referred as the image feature. The text feature from the CLIP text encoder has shape $(C, 512)$ where C is the number of categories. We then apply matrix multiplication between the two to align pixel-level features and a text query feature and perform semantic segmentation using LSeg spatial regularization blocks.

SAM Feature Following the image encoding details in SAM [20], we use an MAE [15] pre-trained ViT-H/16 [8] with 14×14 windowed attention and four equally-spaced global attention blocks. The SAM encoder first obtains the image resolution of 1024×1024 by resizing the image and padding the shorter side. The resolution is then $16 \times$ down-scaled to 64×64 . Since only a portion of the 64×64 feature map contains semantic information due to the padding operation, we crop out one side of the feature map correspond-

ing to the longer side of the original image. Specifically, suppose the original image has the resolution of $H \times W$ where $W > H$, we crop the 64×64 feature map from SAM encoder so that the new feature resolution becomes $64W/H \times 64$ with the feature dimension of 256 corresponding to the output dimension of the SAM encoder. In practice, we use the Speed-up Module with the rendered feature $dim = 128$.

To obtain the results of promptable or promptless segmentation during the inference, we perform the padding operation on the rendered feature to convert from $64W/H \times 64$ back to 64×64 so that the SAM decoder receive the equivalent semantic information as from the original SAM encoder.

Feature visualization As shown in Fig. D, similar to [42], we use `sklearn.decomposition.PCA` [33] in scikit-learn package for feature visualization. We set the number of PCA components to 3 corresponding to RGB channels and calculate the PCA mean by sampling every third element along $h \times w$ vectors, each with feature dimension of either 512 (for LSeg) or 256 (for SAM). The feature map is transformed using the PCA components and mean. This involves centering the features with PCA mean and then projecting them onto PCA components. We then normalize the transformed feature based on the minimum and maximum values with the outliers removed to standardize the feature values into a consistent range so that it can be effectively visualized, typically as an image. We visualize both LSeg and SAM features of scenes from the LLFF dataset [29] from different views. The feature map from LSeg encoder has size 360×480 with dimension 512 (see in the second column of Fig. D). However, as mentioned before, the feature map directly obtained from SAM encoder

Dimension	8	16	32	64	128	256
PSNR \uparrow	36.7061	36.9145	36.8793	36.9208	36.7815	36.9139
SSIM \uparrow	0.9697	0.9706	0.9700	0.9701	0.9700	0.9709
LPIPS \downarrow	0.0234	0.0229	0.0229	0.0228	0.0228	0.0230
FPS \uparrow	64.7	53.9	52.7	32.8	24.2	8.3

Table C. **Evaluation of Image Quality Metrics Across Different Dimensions of SAM feature.** This table presents the PSNR, SSIM, LPIPS, and FPS values corresponding to each dimension level with SAM feature.

contains a padding region (see the red areas on the bottom of the feature maps in the third column), we crop out the region on the feature map as the ground-truth feature before feature distillation (see in the last column). It is worth noting that features from LSeg models mainly capture semantic information by delineating coarse-grained boundaries, while features from SAM models show instance-level information and even fine-grained details in different parts of an object. The capability of teacher encoders determines characteristics of the feature map, thereby influencing the upper limit of the performance of the rendered features on downstream tasks.

D. Replica Dataset Experiment

Following the same selection in [42], we experiment on 4 scenes from the Replica dataset [43]: room_0, room_1, office_3, and office_4. For each scene, 80 images are captured along a randomly chosen trajectory, and every 8th image starting from the third is selected. We trained 5,000 iterations on each scene with LSeg serving as the foundational model for this experiment. We manually re-label some pixels with semantically close labels such as ‘‘rugs’’ and ‘‘floor’’. This preprocess step follows the same method in the NeRF-DFD [21]. The model is trained on the training images and was subsequently evaluated on a set of 10 test images. We test pixel-wise mean intersection-over-union and accuracy on the manually relabeled test images and we use $class = 7$ for the mIoU metric. For room_1, the last 2 test images are excluded from the results since these images do not have 7 classes in the image.

E. Editing Algorithm and Details

The editing procedure takes advantage of the 3D Gaussians so that the model is able to render a novel view image edited with a specific editing operation. As illustrated in Fig. E, starting from a set of 3D Gaussians, i.e. $\mathcal{X} = \{x_1, \dots, x_N\}$ where each x_i is a 3D Gaussian represented by (f_i, α_i, c_i) where $f_i \in \mathbb{R}^{512}$, $\alpha_i \in \mathbb{R}$ and $c_i \in \mathbb{R}^3$ are the semantic feature, color and opacity, respectively. Guided by language, the edit algorithm takes an input text which is a list of object categories, e.g. ‘apple, banana, others’. We leverage

the CLIP’s ViT-B/32 text encoder for text encoding to obtain the text feature $\{t_1, \dots, t_C\}$ where $t_i \in \mathbb{R}^{512}$ and C is the number of categories. We then calculate the inner product of the text feature and semantic feature followed by a softmax function to obtain the semantic scores for each 3D Gaussian, represented by a C -dimensional vector, i.e. $scores \in \mathbb{R}^{N \times C}$. Queried by the text label l , e.g. ‘apple’ or a list of objects to be edited, e.g. ‘apple, banana’, one can either choose to apply hard selection or soft selection to perform edit operation specifically on the target region:

Soft selection: Based on the category selected by the query label $l \in \{1, 2, \dots, C\}$ (or $l \subseteq \{1, 2, \dots, C\}$ if l is a list of categories), we target on the corresponding column of the score matrix, i.e. $score_l = [s_{1l}, s_{2l}, \dots, s_{Nl}]^T$ and apply binary thresholding on this column score vector, i.e. for any i such that $s_{il} \geq th$, we set the position i to 1 representing being selected; otherwise the position i is set to 0 representing not being selected. Then all the positions i such that $s_{il} = 1$ compose a target region to be edited. Intuitively, we mask out all the 3D Gaussians that are not selected and use those selected 3D Gaussians to update the color set $\{c_i\}_{i=1}^N$ and opacity $\{\alpha_i\}_{i=1}^N$.

Hard selection: We apply $argmax$ function to the score matrix to select the category corresponding to the highest score for each Gaussian to obtain a filtered category vector $categories = [c_1, c_2, \dots, c_N]^T$ where $c_i = argmax\{s_{i1}, s_{i2}, \dots, s_{iC}\}$. Then we filter based on the query label l : for any i such that $c_i = l$ (or $c_i \in l$ if l represents a list of target objects), we set the position i to 1 representing being selected; otherwise the position i is set to 0 representing not being selected. Similarly, we select the target region to be edited by preserving only the region positions of which the highest score category is aligned with the query label.

Hybrid selection: Since soft selection applies thresholding only based on column vector of score matrix corresponding to label l which may potentially cause incorrect selection when the dominant score exists in other columns while hard selection merely selects the highest score without any tunable threshold value. Therefore, we propose a hybrid selection method by combining both hard and soft selection to alleviate the effect of incorrectly selecting the

category while making the selection tunable to adapt to different scenarios. Specifically, we combine the selection Gaussian masks of the two methods and apply bitwise OR operation between the two masks to obtain the final selected region.

We then update the opacity and color based on the selected edit region and a specific edit operation. We demonstrate the details of three examples: extraction, deletion and appearance modification:

(a) Extraction: for any $i \in \{1, \dots, N\}$, if i is selected, the opacity remains to be α_i ; otherwise the opacity is set to 0.

(b) Deletion: for any $i \in \{1, \dots, N\}$, if i is selected, the opacity is set to 0; otherwise the opacity remains to be α_i . Specifically for deletion operation, we apply the hybrid selection method to select the target edit region to reduce the effect of incomplete deletion caused by the absence of target pixels.

(c) Appearance modification: for any $i \in \{1, \dots, N\}$, if i is selected, the color c is updated to be $appearance_func(c_i)$ where $appearance_func(\cdot)$ represents any appearance modification function, such as changing the green leaves to red leaves, etc.

F. Ablation Studies

We study the effect of different rendered feature dimensions using our Speed-up Module. In Tab. A, we report the performance of the semantic segmentation on Replica dataset using LSeg feature. The result shows that both rendered feature $dim = 256$ and $dim = 128$ can achieve the best performance on accuracy, and $dim = 256$ is slightly better on mIoU. However, $dim = 128$ is $\times 2.4$ faster than $dim = 256$ on training. We also report the quantitative results of novel view synthesis in Tab. B, which shows that $dim = 128$ is the best. Therefore, we choose $dim = 128$ for our Replica dataset experiment in practice. In addition, we show the performance and speed (FPS) of novel view synthesis with different dimensions of SAM feature in Tab. C.

Furthermore, Fig. A and Fig. B substantiate that our Speed-up Module not only avoids compromising performance but, in fact, resulting in time savings.

G. Failure Cases

The proposed method indeed has limitations reflected on some failure cases. In Fig. C, we showcase failure cases for scenes that are more challenging and complex. In Fig. C (a), the point-prompted segmentation mask is not perfect with a coarse boundary and small holes. This is caused by low feature quality from SAM distillation, rather than Gaussian representation. Since the boundary of the car is hard to depicted and there are multiple similar objects close to each other (multiple adjacent cars), making the scene complex. As a result, achieving a smooth and accurate mask boundary of the car becomes challenging, which could be counted as a limitation. In Fig. C (b), given the text prompt

“Delete the cup”, although succeeding in locating the target object, the model fails to remove the cup comprehensively. The reason behind is that in some complex scenes including various objects with multiple sizes, the 3D Gaussians corresponding to the tiny objects with sophisticated details are hard to accurately selected by the “Gaussian mask”. As a result, a clean deletion is hard to perform on target object.

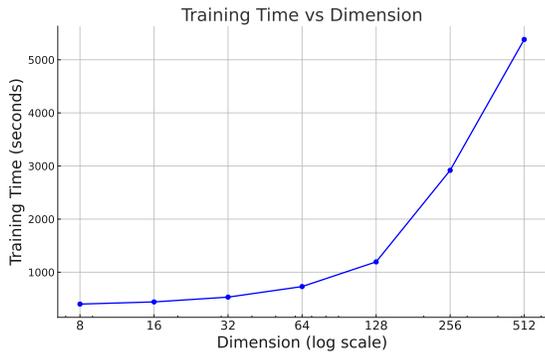


Figure A. **Training Time vs Speed-up Module Dimension** We test the training time required with different input dimension of speed-up module. In this Figure, we show that the training time can be significantly reduced with our speed-up module.

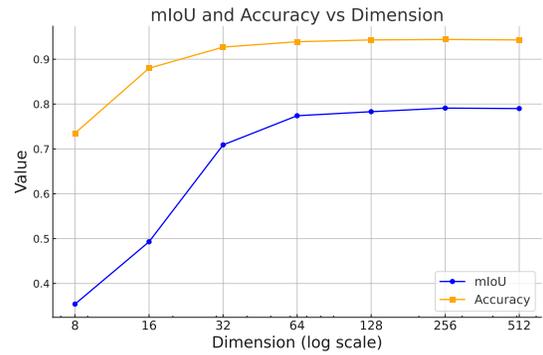


Figure B. **mIoU and Accuracy vs Dimension** In this graph, we show 2D metrics with respect to different input dimensions of speed-up module. With our speed-up module and proper input dimension, the 2D metrics are not compromised.

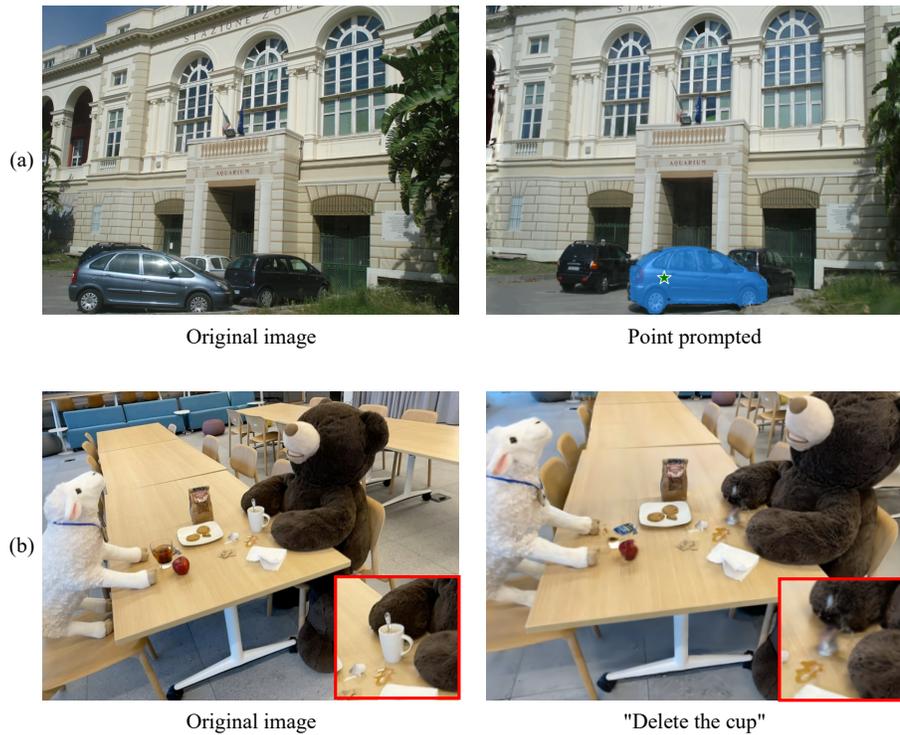
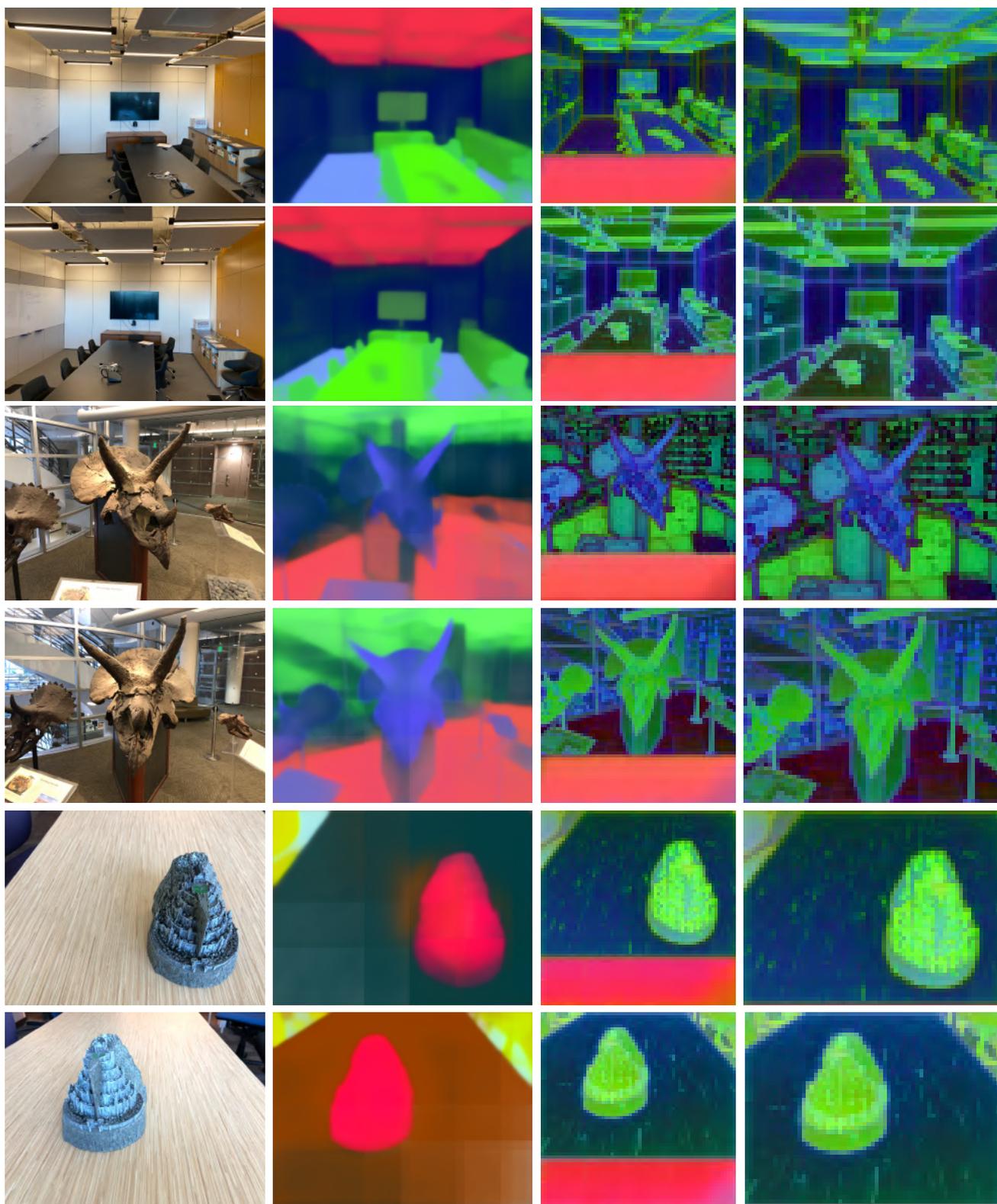


Figure C. **Failure cases in complex and challenging situations** (a) The point-prompted segmentation mask, contains flaws in the form of a coarse boundary and small holes, resulting from low-quality features. (b) The model fails to delete tiny sophisticated objects thoroughly in a complex scene in language-guided editing.



(a) Original image

(b) LSeg feature

(c) SAM feature

(d) SAM cropped feature

Figure D. **Feature visualization on different scenes from LLFF dataset [29] from LSeg and SAM encoders.** Note that SAM features in column (d) is obtained by cropping the padding region. We resize the cropped feature for better visualization.

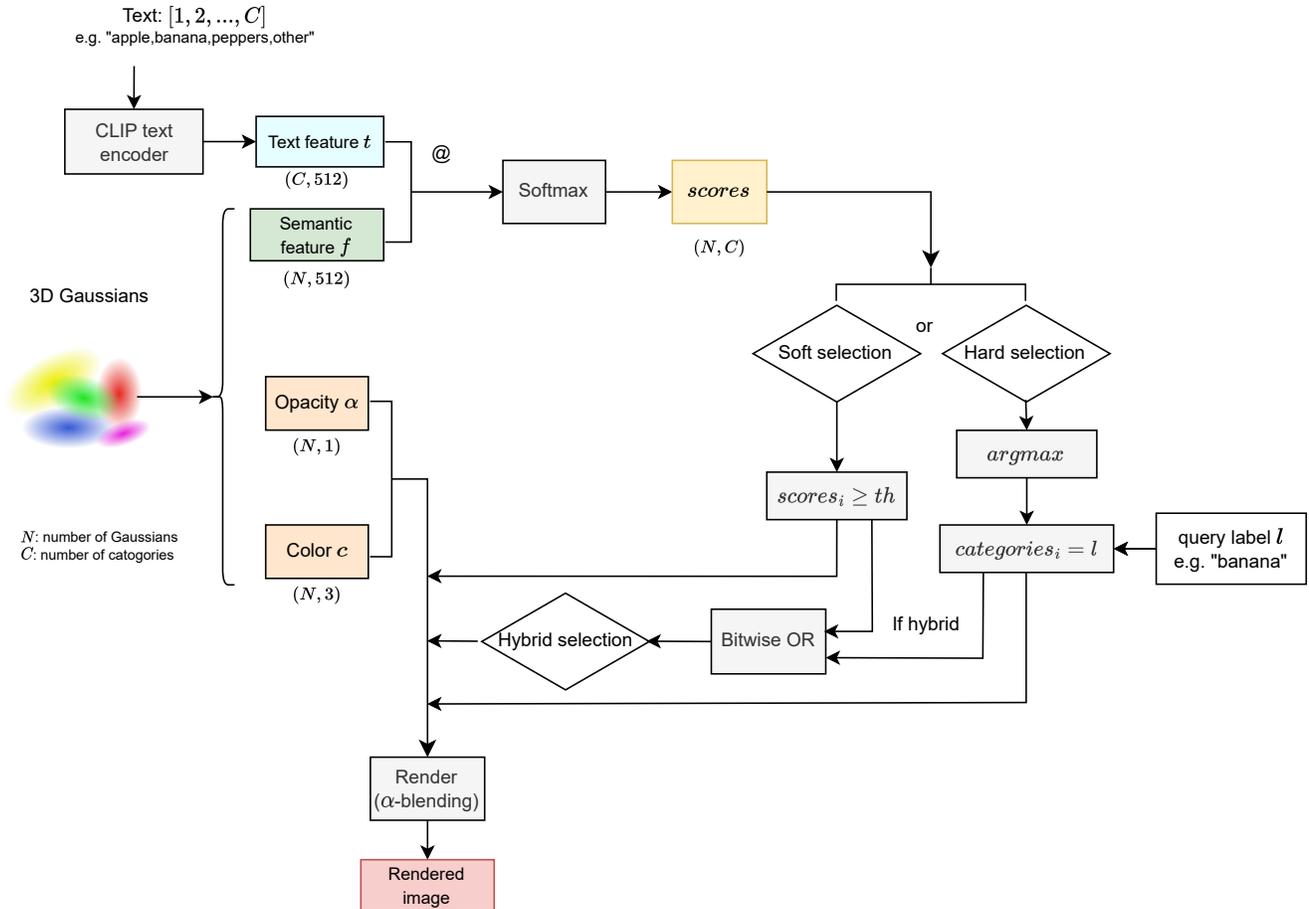


Figure E. **Language-guided editing procedure using 3D Gaussians.** We calculate the inner product between the semantic feature and the text feature from CLIP encoder followed by a softmax to obtain a score matrix and query the feature field to apply editing on target regions (obtained from soft selection / hard selection / hybrid selection) by updating opacity and color from 3D Gaussians before rendering.