# Gradient-Weighted Feature Back-Projection: A Fast Alternative to Feature Distillation in 3D Gaussian Splatting

Joji Joseph
Indian Institute of Science
jojijoseph@iisc.ac.in

Bharadwaj Amrutur
Indian Institute of Science
amrutur@iisc.ac.in

Shalabh Bhatnagar
Indian Institute of Science
shalabh@iisc.ac.in

## Abstract

*We introduce a training-free method for feature field rendering in Gaussian splatting. Our approach back-projects 2D features into pre-trained 3D Gaussians, using a weighted sum based on each Gaussian's influence in the final rendering. While most training-based feature field rendering methods excel at 2D segmentation but perform poorly at 3D segmentation without post-processing, our method achieves high-quality results in both 2D and 3D segmentation. Experimental results demonstrate that our approach is fast, scalable, and offers performance comparable to training-based methods. Project page:* [https://jojijoseph.github.io/3dgs-backprojection/](https://jojijoseph.github.io/3dgs-backprojection/)

## 1. Introduction

3D Gaussian Splatting (3DGS) [3] is a novel-view synthesis technique that uses 3D Gaussians as rendering primitives, each defined by its mean position and variance. Additionally, these Gaussians carry payloads such as opacity and anisotropic color, enabling photorealistic rendering.

Generating intermediate feature maps directly, rather than passing RGB renderings to a foundation model, can be advantageous for tasks like segmentation. However, training Gaussian splatting models to render feature maps directly is computationally intensive, especially given the high dimensionality of feature space.

Moreover, segmentation queries are more effective when performed on rendered features than on features assigned to each Gaussian, as only the surface Gaussians typically align with the feature map. Since rendered features result in a weighted sum, Gaussians deeper within the scene may not accurately correspond to features in the feature map, complicating accurate 3D segmentation.

Given these drawbacks in training time and segmentation accuracy, alternative methods are essential. One approach is to aggregate 2D features into Gaussians that influence the regions where the 2D features are projected. In other words,

the Gaussian features are computed as a weighted sum of the 2D features, with weights proportional to the influence of each Gaussian on the corresponding region in the 2D feature map.

Our main contributions are as follows:

- We introduce a training-free approach that projects 2D features onto 3D Gaussians, providing a fast and scalable alternative to traditional feature field distillation.
- We demonstrate that our method achieves comparable or superior performance to feature field distillation methods that rely on extensive training, particularly in producing clean 3D segmentations.
- Our method effectively supports both rendered feature space and 3D feature space, enabling seamless querying for downstream applications such as 3D object manipulation and real-time scene understanding.

## 2. Related Works

Neural Radiance Fields (NeRF) [12] has emerged as a leading approach for synthesizing novel views from sparse image inputs. By leveraging a neural network, NeRF captures a volumetric representation of a scene, modeling the color and density at every 3D location to generate highly realistic views from arbitrary perspectives. However, the implicit nature of NeRF's scene representation poses challenges for tasks like object modification or rearrangement, as such operations typically demand retraining the entire network.

In contrast, 3D Gaussian Splatting [8] provides an explicit representation of 3D scenes, using 3D Gaussians as the primary rendering primitives. Attributes like color, opacity, and orientation characterize these Gaussians. This explicit structure allows for direct manipulation of Gaussians and their associated parameters, enabling efficient object rearrangement and editing. Such flexibility makes 3D Gaussian Splatting well-suited for interactive applications, including object editing, augmented reality, digital twins, and robotics.

An intuitive progression from radiance field rendering is feature field rendering, which incorporates additional fea-

ture embeddings to enrich the representation. Recent studies, such as [15, 17, 20], have advanced this concept for tasks like segmentation and semantic querying.

Feature-3DGS [20] focuses on training high-dimensional feature embeddings, while LangSplat [15] prioritizes compressed, low-dimensional features. Both methods demonstrate strong performance in 2D segmentation of rendered outputs but face significant challenges with 3D object segmentation. Feature-3DGS attempts 3D segmentation by matching language embeddings with Gaussian feature embeddings. However, this approach often falls short because the features of individual Gaussians do not directly map to the final rendered feature, which results from the weighted sum of contributions from multiple Gaussians (see Equation 2). This inherent mismatch hinders reliable 3D segmentation. In contrast, our method overcomes this limitation by directly leveraging gradient information, enabling more accurate and effective 3D segmentation.

Another class of methods tackles object segmentation from visual prompts such as points and masks[1, 6, 7, 16]. In SAGD [6], a binary voting system is employed, while methods in [16] and [7] use similar influence-based voting approaches. FlashSplat [16] formulates segmentation as an optimization problem, whereas [7] leverages inference-time gradient backpropagation. Our approach also utilizes inference-time backpropagation for feature back-projection, enabling robust feature representation across 3D space.

## 3. Method

In this section, we present the feature back-projection equation and describe four direct use cases — 3D segmentation, affordance transfer, and identity encoding—that do not require post-processing other than similarity search. Notably, we perform these use cases directly in 3D space rather than in 2D image space, as is common in similar works.

### 3.1. Feature Back-Projection

Consider the color $C$ of a pixel at $(x, y)$ in a 3DGS rendering,

$$C(x, y) = \sum_{n \leq N} c_n \alpha_n(x, y) \prod_{m < n} (1 - \alpha_m(x, y)) \quad (1)$$

$$= \sum_{n \leq N} c_n \alpha_n(x, y) T_n(x, y) \quad (2)$$

Where $N$ is the total number of Gaussians, each indexed by its sorted position, $c_n$ is the color associated with the $n$th Gaussian, $\alpha_n(x, y)$ is the opacity of the $n$th Gaussian at $(x, y)$ adjusted with exponential falloff, and $T_n = \prod_{m < n} (1 - \alpha_m(x, y))$ is the transmittance of $n$th Gaussian at $(x, y)$.

From this equation, it's clear that the rendered color is a weighted sum of colors of individual Gaussians. Moreover, the weight is opacity multiplied by transmittance.

Taking the derivative with respect to color of $k$th Gaussian $c_k$,

$$\frac{\partial C(x, y)}{\partial c_k} = \alpha_k(x, y) T_k(x, y) \quad (3)$$

This gradient is equivalent to the weight of each Gaussian in a pixel rendering. This insight enables us to leverage inference-time gradient backpropagation to compute feature back-projections based on the Gaussian's influence efficiently.

Given this gradient-based weighting, we can define the feature back-projection equation as follows:

$$\mathbf{f}_k = \frac{\sum_{(x,y,n)} \mathbf{F}_{2D}(x, y, n) \alpha_k(x, y, n) T_k(x, y, n)}{\sum_{(x,y,n)} \alpha_k(x, y, n) T_k(x, y, n)} \quad (4)$$

We call this equation the expected feature back-projection equation or simply the back-projection equation.

Where $\mathbf{f}_k$ is the feature of $k$th Gaussian. $\mathbf{F}_{2D}(x, y, n)$ is the feature at $(x, y)$ in the $n$th viewpoint. This formulation allows for efficient aggregation of features across viewpoints, weighted by opacity and transmittance, resulting in an accurate feature back-projection that reflects the Gaussian's contribution to the final rendered scene.

When we remove the denominator from this equation, it becomes accumulated back-projection.

$$\mathbf{f}_k = \sum_{(x,y,n)} \mathbf{F}_{2D}(x, y, n) \alpha_k(x, y, n) T_k(x, y, n) \quad (5)$$

If we replace the feature $\mathbf{F}_{2D}(x, y, n)$ with a function indicating binary function that indicates the presence of a mask this simply becomes the vote using masked gradients as described in [7].

If we do Euclidian normalization on $\mathbf{f}_k$ after calculation, both equations 4, 5 become equivalent.

### 3.2. 3D Segmentation

Once the feature back-projection is complete, we obtain a set of feature vectors of dimension $D$, $\{\mathbf{f}_k\} \in \mathbb{R}^D$ in the projected feature space. Let $\mathbf{q} \in \mathbb{R}^D$ represent the embedding extracted from the language query (or a query from another modality). The scalar value $\text{sim}(\mathbf{f}_k, \mathbf{q})$, where sim denotes a similarity function (typically cosine similarity), measures the relevance of each Gaussian $g_k$ to the query.

To perform segmentation, we apply a threshold on $\text{sim}(\mathbf{f}_k, \mathbf{q})$, allowing us to isolate Gaussians that correspond closely to the specified query. Formally, the set of segmented Gaussians $G$ can be defined as:

$$G = \{g_k \mid \text{sim}(\mathbf{f}_k, \mathbf{q}) > \theta\}, \quad (6)$$

where $\theta$ is a user-defined threshold.

Segmentation can also be performed in 2D. The equation 7 represents the 2D mask by querying over rendered 2D features.

$$M(x, y) = \begin{cases} 1, & \text{if } \text{sim}(\mathbf{f}_k, \mathbf{q}) > \theta \\ 0, & \text{otherwise} \end{cases} \qquad (7)$$

Here, $M(x, y)$ is a binary mask that identifies regions in 2D space where the similarity to the query exceeds the threshold.

### 3.3. Affordance Transfer

Affordance transfer[5] is the process of transferring annotated regions from source images to a target object. Although the source and target objects may differ in appearance, they belong to the same category, enabling the transfer of meaningful information across instances. This technique is a particular case of few-shot segmentation, commonly applied in robotics, where annotations of irregularly shaped regions are transferred to a target object. This transfer allows a robot to manipulate the target object using knowledge from the annotated regions.

Our approach follows a straightforward pipeline. First, we transfer DINOv2[4, 14] features to the target scenes via feature back-projection. Then, we identify the closest matching source annotation for each Gaussian in the target scene by performing a k-nearest neighbors (kNN) classification on the feature space. Here, the kNN classification finds source annotations most similar to the Gaussian features, ensuring that the transferred affordances correspond closely to the target object's structure.

### 3.4. Identity Encoding

When annotations are available for different objects within a scene, we can encode each object with a unique vector and back-project these vectors into the 3D scene. Even if annotations are missing for some views, this method can generalize across the entire 3D scene given sufficient annotated views.

We implement two types of identity encoding, each with unique strengths:

**Orthogonal Encoding**: In this approach, we assign procedurally generated orthogonal vectors to represent each object uniquely, where each vector direction is associated with a distinct object. The orthogonality of these vectors ensures that each object has a distinct, non-overlapping identity in the feature space, maximizing separation and simplifying the segmentation task. This back-projection operation can be viewed as an extension of segmentation with masked gradients [7], enabling the simultaneous encoding of multiple classes within a single scene. While effective, orthogonal encoding becomes challenging when representing a large number of objects, as the dimensionality of the feature space limits orthogonal vectors.

**Contrastive Encoding**: For cases where we need to encode many objects, we use non-orthogonal vectors made sufficiently distinct through contrastive learning. Contrastive learning maximizes the distance between feature vectors of different objects by pulling representations of different objects apart and bringing representations of the same object closer together. This method allows us to encode more objects than possible with strictly orthogonal vectors while still achieving clear separation between object identities in feature space. Although the separation is not as strict as with orthogonal encoding, contrastive learning provides a scalable and effective alternative, maintaining discriminative power even with far more objects than the embedding dimension.

The loss function for training the Identity Encoder-Decoder model is defined as:

$$\mathcal{L} = \mathcal{L}_{\text{classification}} + \mathcal{L}_{\text{orthogonality}}, \qquad (8)$$
$$\mathcal{L}_{\text{classification}} = \text{CrossEntropy}(\hat{y}, y), \qquad (9)$$
$$\mathcal{L}_{\text{orthogonality}} = \|\mathbf{E}\mathbf{E}^{\top} - \mathbf{I}\|_F. \qquad (10)$$

Where $\hat{y}$ is the model's predicted class probability and $y$ is the true class label. $\mathbf{E}$ is an embedding matrix of size $N \times D$, where $N$ is the number of classes and $D$ is the embedding dimension. $\mathbf{I}$ is the identity matrix and $\| \cdot \|_F$ denotes the Frobenius norm.

## 4. Experiments and Results

### 4.1. Setup

We use gsplat[19] as our rasterizer. All experiments are performed on NIVIDA A6000 GPU.

We use LSeg[10], DINOv2[4, 14] features for our experiments. We also show we can encode objects using procedurally generated features as in orthogonal encoding and latent features generated by contrastive encoding.

There are cases where the numerator in the back-projection equation(4) is zero since there is no influence of Gaussian in the screen. In those cases, we can prune out the Gaussians with zero numerators before proceeding with inference.

### 4.2. 3D Segmentation

In this section, we use LSeg[10] features for our segmentation. We compare it qualitatively with Feature 3DGS. We use a custom implementation of Feature 3DGS based on gsplat [19] to prevent excessive RAM usage. We train Feature 3DGS versions for 7000 iterations. The training process takes around 20-30 minutes on our system, while the same number of iterations of Vannila 3DGS, which doesn't

support feature field training, takes only 2 minutes 30 seconds. That is 10x faster compared to using Feature 3DGS.

After back-projection, we do segmentation as follows,

Let $\mathbf{Q} = \{\mathbf{q}_i\}$ be the set of prompts. Let $\mathbf{q}_0$ be the category to segment, $\mathbf{q}_i, i \neq 0$ are negative prompts. We use the last prompt as 'other'.

The 3D mask is found by taking cosine similarity with Gaussian with both positive and negative prompts. Then, take the Gaussians, where the category prompt gives the highest cosine similarity as the 3D mask. Additionally, we can use a threshold over the cosine similarity with the category prompt.

The qualitative 2D segmentation results are given in figure 1, and 3D segmentation results are given in figure 2. For 2D segmentation, both methods produce similar results. But for the 3D segmentation, our method gives better results.

The feature back-projection is completed in a single pass through the training views, with the entire process—including ground truth feature map generation—taking an average of 2-3 minutes. During inference, our method achieves object segmentation in just 30 ms, including both text encoding and similarity calculation. This represents a 900x speedup compared to segmentation using masked gradients [7], a method designed for clean segmentation, when considering the full pipeline including mask generation. In our approach, the majority of the latency arises from text encoding, while similarity calculation is highly efficient, requiring less than 1 ms.

### 4.3. Affordance Transfer

We transfer DINOv2 features to each scene of the [13] dataset. Then use source images and annotations from [7] and transfer it directly to the Gaussians of target scene. We label the method from [7] as 2D-2D-3D transfer because source annotations are transferred to 2D target frames before applying to 3D and our method as 2D-3D for contrasting.
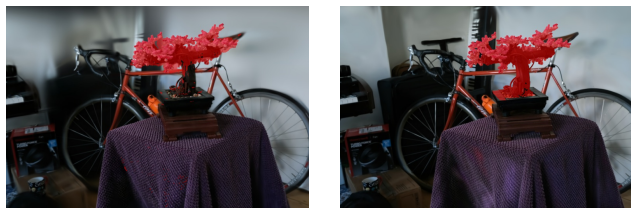
Table 1. Table showing comparison affordance transfer with masked gradients and our method. The label 2D-3D indicate our method.

| Scene | mIoU ↑ | | Recall ↑ | | Time Taken (s) ↓ | |
|---|---|---|---|---|---|---|
| | 2D-2D-3D | 2D-3D | 2D-2D-3D | 2D-3D | 2D-2D-3D | 2D-3D |
| 1 | **47.87** | 42.80 | **67.77** | 67.11 | 293.88 | **5.22** |
| 2 | **55.63** | 53.28 | 81.07 | **82.55** | 317.12 | **8.03** |
| 3 | **60.50** | 57.82 | **86.95** | 86.68 | 142.79 | **7.58** |
| Mean | **54.67** | 51.30 | 78.60 | **78.78** | 251.26 | **6.94** |

### 4.4. Identity Encoding

**Orthogonal Encoding**: For simplicity we use one-hot encoding for identity encoding. That is one element is 1 and



**2D Segmentation (F3DGS)**     **2D Segmentation (Ours)**

Prompt: Plant, Negative Prompt: Other

Prompt: Table, Negative Prompt: Vase, Other

Prompt: Plant, Negative Prompt: Vase, Other

Prompt: Truck, Negative Prompt: Other

Figure 1. Comparison of 2D segmentation in Feature 3DGS (F3DGS) and our method. Under each pair of image corresponding positive and negative prompts are given.

rest of them are 0. The qualitative results are shown in figure 4. The scene is taken from 3D-OVS dataset[11].

**Contrastive Encoding**: We use LERF-Mask dataset introduced for the identity encoding method Gaussian Grouping[18] for quantitative evaluation. Here we use 16 dimension embeddings for each group. Total number of groups are around 200.

We first train a classifier to predict the group over the embeddings. We make sure to use a contrastive loss to make embeddings far apart from each other. Then back-project this embeddings to each Gaussian. We use the classifier to predict the groups each rendered pixel belongs.

We follow the evaluation protocol from [18] to evaluate our method. See table 2 for the quantitative evaluation.

4

| Extraction (F3DGS) | Extraction (Ours) | Deletion (F3DGS) | Deletion (Ours) |
|---|---|---|---|



Prompt: Plant, Negative Prompt: Other

Prompt: Table, Negative Prompt: Vase, Other

Prompt: Plant, Negative Prompt: Vase, Other
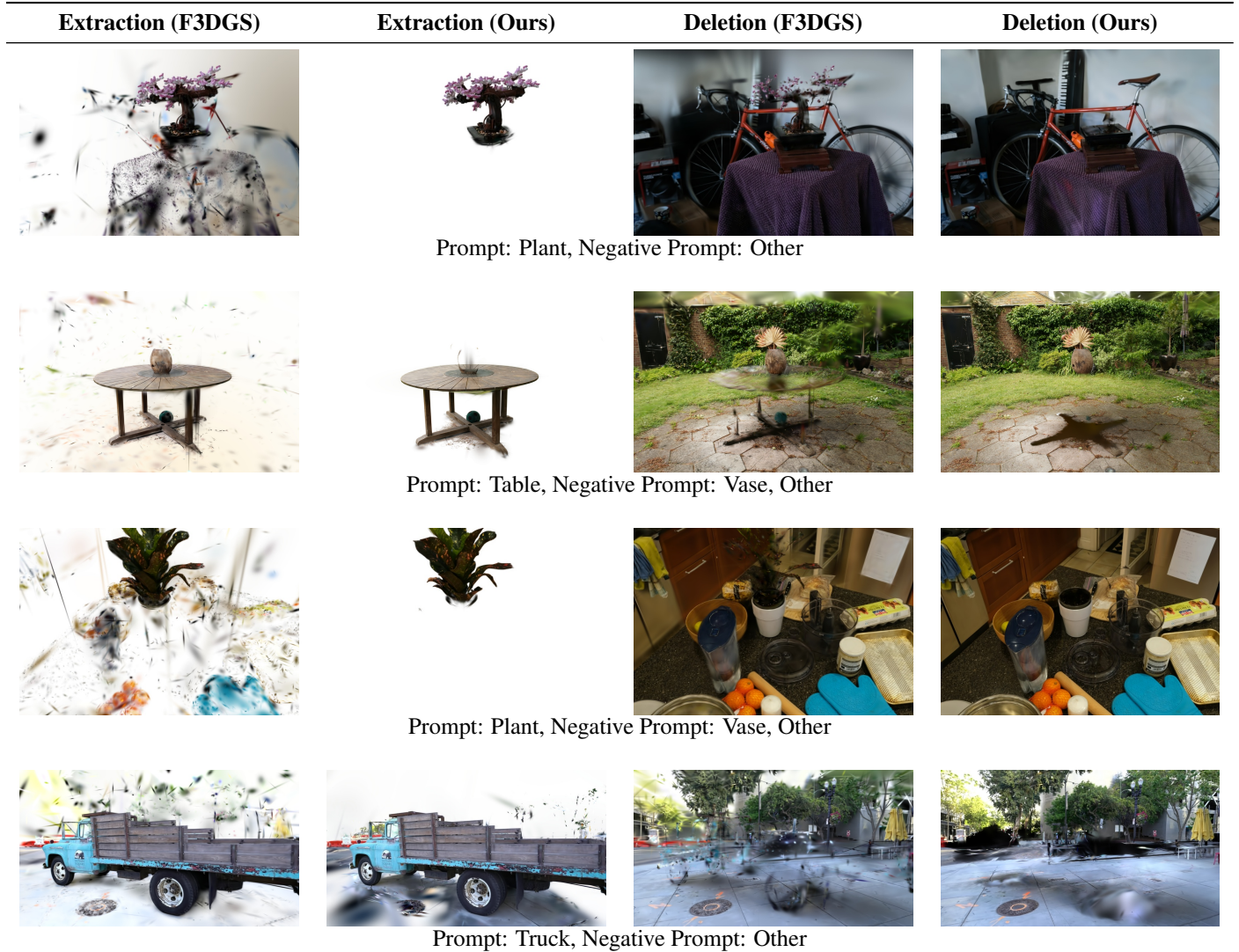
Prompt: Truck, Negative Prompt: Other

Figure 2. Comparison of 3D segmentation in Feature 3DGS (F3DGS) and our method. Under each pair of image corresponding positive and negative prompts are given. Clearly our method produces less outliers.

Note that we evaluate using 2D mIoU. Surprisingly we are able to get similar performance as that of Gaussian Grouping which trains identity encoding along with each scene. Our method takes only 10 seconds for training the classifier and 10 seconds to transfer the identity encoding.

## 5. Discussion and Conclusion

We have introduced a novel, training-free, efficient, scalable alternative to feature field distillation in Gaussian splatting. Our method achieves fast, clean segmentation by directly querying the features associated with each Gaussian.

Our approach aggregates features from all available training views in a single pass, mitigating any inconsisten-

Table 2. Comparison of mIoU scores for different models across Figurines, Ramen, and Teatime datasets. Bold indicates best performance and underline indicates second best. Our method is comparable to the state of the art Gaussian Grouping[18]

. The readings other that ours are from [18].

| Model | Figurines | Ramen | Teatime | Mean |
|---|---|---|---|---|
| DEVA [3] | 46.2 | 56.8 | 54.3 | 52.4 |
| LERF [9] | 33.5 | 28.3 | 49.7 | 37.2 |
| SA3D [2] | 24.9 | 7.4 | 42.5 | 24.9 |
| LangSplat [15] | 52.8 | 50.4 | 69.5 | 57.6 |
| Gaussian Grouping [18] | 69.7 | **77.0** | 71.7 | 72.8 |
| Ours | **73.5** | 72.7 | **74.1** | **73.4** |

(a) Source Images with anno-    (b) Target scene with transferred an-
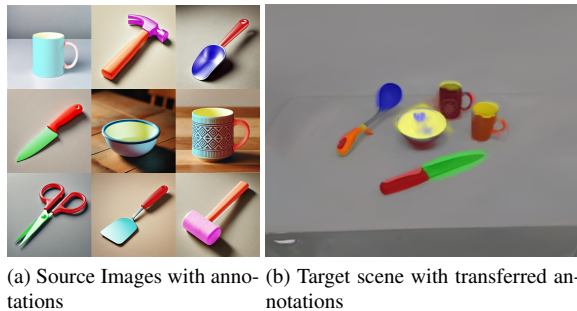tations                          notations

Figure 3. Qualitative result of affordance transfer. The example images are different but the same category as that of target scene.



(a) Original View          (b) Rendering with grouping
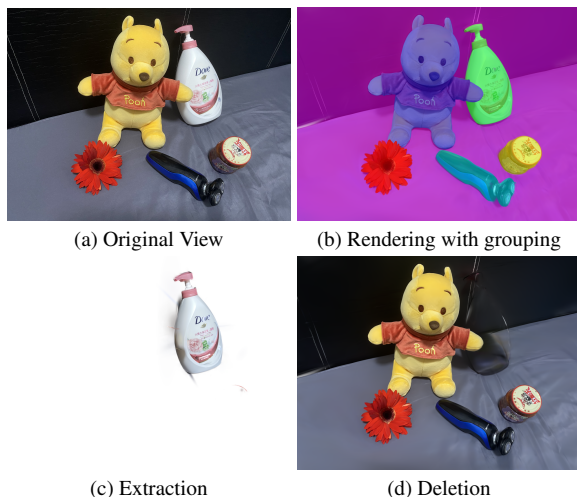
(c) Extraction             (d) Deletion

Figure 4. Sample result of orthogonal identity encoding. We use only about 5 frames with ground truth. The scene is trained with 30 imags. We are able to extract the occluded object as well as delete it. Note that the results are shown without post processing to remove outliers.

cies in individual 2D feature maps through an averaging effect. Nevertheless, minor imperfections may still propagate, though they have minimal impact on overall performance.

Since our method does not update Gaussian parameters, it does not benefit from the regularization effect inherent in traditional feature field methods. However, we found no significant drawbacks when compared to feature field methods—in fact, our approach is considerably faster for generating feature embeddings for Gaussians and produces superior qualitative results.

## References

[1] Jiazhong Cen, Jiemin Fang, Chen Yang, Lingxi Xie, Xiaopeng Zhang, Wei Shen, and Qi Tian. Segment any 3d gaussians. *arXiv preprint arXiv:2312.00860*, 2023. 2

[2] Jiazhong Cen, Zanwei Zhou, Jiemin Fang, Chen Yang, Wei Shen, Lingxi Xie, Dongsheng Jiang, Xiaopeng Zhang, and Qi Tian. Segment anything in 3d with nerfs. In *NeurIPS*, 2023. 5

[3] Ho Kei Cheng, Seoung Wug Oh, Brian Price, Alexander Schwing, and Joon-Young Lee. Tracking anything with decoupled video segmentation. In *ICCV*, 2023. 5

[4] Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers, 2023. 3

[5] Denis Hadjivelichkov, Sicelukwanda Zwane, Marc Deisenroth, Lourdes Agapito, and Dimitrios Kanoulas. One-Shot Transfer of Affordance Regions? AffCorrs! In *Proceedings of The 6th Conference on Robot Learning (CoRL)*, pages 550–560, 2023. 3

[6] Xu Hu, Yuxi Wang, Lue Fan, Junsong Fan, Junran Peng, Zhen Lei, Qing Li, and Zhaoxiang Zhang. Sagd: Boundary-enhanced segment anything in 3d gaussian via gaussian decomposition, 2024. 2

[7] Joji Joseph, Bharadwaj Amrutur, and Shalabh Bhatnagar. Gradient-driven 3d segmentation and affordance transfer in gaussian splatting from 2d masks. *arXiv preprint arXiv:2409.11681*, 2024. 2, 3, 4

[8] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42 (4), 2023. 1

[9] Justin* Kerr, Chung Min* Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lerf: Language embedded radiance fields. In *International Conference on Computer Vision (ICCV)*, 2023. 5

[10] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and Rene Ranftl. Language-driven semantic segmentation. In *International Conference on Learning Representations*, 2022. 3

[11] 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. *arXiv preprint arXiv:2305.14093*, 2023. 4

[12] 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. In *ECCV*, 2020. 1

[13] Austin Myers, Ching L. Teo, Cornelia Fermüller, and Yiannis Aloimonos. Affordance detection of tool parts from geometric features. In *ICRA*, 2015. 4

[14] Maxime Oquab, Timothée Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2023. 3

[15] Minghan Qin, Wanhua Li, Jiawei Zhou, Haoqian Wang, and Hanspeter Pfister. Langsplat: 3d language gaussian splatting. *arXiv preprint arXiv:2312.16084*, 2023. 2, 5

[16] Qiuhong Shen, Xingyi Yang, and Xinchao Wang. Flashsplat: 2d to 3d gaussian splatting segmentation solved optimally. *European Conference of Computer Vision*, 2024. 2

[17] Jin-Chuan Shi, Miao Wang, Hao-Bin Duan, and Shao-Hua Guan. Language embedded 3d gaussians for open-vocabulary scene understanding. *arXiv preprint arXiv:2311.18482*, 2023. 2

[18] Mingqiao Ye, Martin Danelljan, Fisher Yu, and Lei Ke. Gaussian grouping: Segment and edit anything in 3d scenes. In *ECCV*, 2024. 4, 5

[19] Vickie Ye, Ruilong Li, Justin Kerr, Matias Turkulainen, Brent Yi, Zhuoyang Pan, Otto Seiskari, Jianbo Ye, Jeffrey Hu, Matthew Tancik, and Angjoo Kanazawa. gsplat: An open-source library for Gaussian splatting. *arXiv preprint arXiv:2409.06765*, 2024. 3

[20] Shijie Zhou, Haoran Chang, Sicheng Jiang, Zhiwen Fan, Zehao Zhu, Dejia Xu, Pradyumna Chari, Suya You, Zhangyang Wang, and Achuta Kadambi. Feature 3dgs: Supercharging 3d gaussian splatting to enable distilled feature fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21676–21685, 2024. 2