# LangSplatV2: High-dimensional 3D Language Gaussian Splatting with 450+ FPS

**Wanhua Li**[1,*]    **Yujie Zhao**[1,2,*]    **Minghan Qin**[3,*]    **Yang Liu**[3]    **Yuanhao Cai**[4]
**Chuang Gan**[5,6]    **Hanspeter Pfister**[1]

[1]Harvard University    [2]University of Chinese Academy of Sciences    [3]Tsinghua University
[4]Johns Hopkins University    [5]MIT-IBM Watson AI Lab    [6]UMass Amherst

## Abstract

In this paper, we introduce LangSplatV2, which achieves high-dimensional feature splatting at 476.2 FPS and 3D open-vocabulary text querying at 384.6 FPS for high-resolution images, providing a $42 \times$ speedup and a $47 \times$ boost over LangSplat respectively, along with improved query accuracy. LangSplat employs Gaussian Splatting to embed 2D CLIP language features into 3D, significantly enhancing speed and learning a precise 3D language field with SAM semantics. Such advancements in 3D language fields are crucial for applications that require language interaction within complex scenes. However, LangSplat does not yet achieve real-time inference performance (8.2 FPS), even with advanced A100 GPUs, severely limiting its broader application. In this paper, we first conduct a detailed time analysis of LangSplat, identifying the heavyweight decoder as the primary speed bottleneck. Our solution, LangSplatV2 assumes that each Gaussian acts as a sparse code within a global dictionary, leading to the learning of a 3D sparse coefficient field that entirely eliminates the need for a heavyweight decoder. By leveraging this sparsity, we further propose an efficient sparse coefficient splatting method with CUDA optimization, rendering high-dimensional feature maps at high quality while incurring only the time cost of splatting an ultra-low-dimensional feature. Our experimental results demonstrate that LangSplatV2 not only achieves better or competitive query accuracy but is also significantly faster. Codes and demos are available at our project page: `https://langsplat-v2.github.io`.

## 1 Introduction

Seamless and intuitive interactions between humans and complex 3D environments [1, 2] are paramount for a wide range of applications, such as augmented reality [3, 4] and intelligent robotics [5, 6, 7]. To achieve such interactions, systems must understand and respond to natural language queries in real time. This capability hinges on advancements in 3D language fields—a technology at the intersection of vision-language models [8, 9] and 3D environmental modeling [10, 11].

LangSplat [12], a pioneering model in this domain, leverages 3D Gaussian Splatting to embed 2D CLIP language features into 3D spaces. This method has significantly accelerated the 3D query time, achieving speeds up to $199 \times$ faster than its predecessors [13], which is critical for applications in scenarios where rapid response is essential. Despite these advancements, LangSplat does not achieve real-time performance, especially at high resolutions, which hinders its widespread application.

Real-time querying is crucial for applications such as real-time navigation [6], interactive gaming [14], and on-the-fly educational tools [15], where delays can disrupt user experience and functionality. To identify the speed bottleneck of LangSplat, we decompose the entire querying process into three

---

stages: rendering, decoding, and post-processing. Then, we provide a detailed time analysis on the LERF dataset with one A100 GPU. Results in Table 1 show that each query takes 122.1 ms for LangSplat. With some simple engineering modifications, the rendering and post-processing time can be significantly reduced, we term this new version of LangSplat as LangSplat*. The performance of LangSplat* revealed that the decoding stage, which is responsible for transforming low-dimensional latent features into high-dimensional CLIP features with a heavy-weight multi-layer perceptron (MLP), is the primary bottleneck. LangSplat introduces the MLP decoder to significantly reduce the training memory and time cost. Consequently, an additional decoding stage with an MLP is inevitable for test-time querying and reducing the dimensionality of high-dimensional features lowers the accuracy of the queries. However, simply removing the decoder and directly training the high-dimensional CLIP feature for each Gaussian dramatically increases the rendering time. As shown in Figure 1, as the feature splatting dimension increases, the rendering time of LangSplat significantly increases. Specifically, rendering a feature with a dimension of 1536 (assuming we render three semantic levels of 512-dimensional CLIP feature fields in parallel) is 15 times slower than rendering a 3-dimensional field on one A100 GPU.

To address the decoding bottleneck in LangSplat, we introduce LangSplatV2, which drastically enhances querying speed without sacrificing querying accuracy. As shown in Figure 1, our LangSplatV2 successfully decouples rendering speed from the dimensionality of rendering features, enabling the rendering of high-dimensional features at the computational cost of splatting an ultra-low-dimensional feature. As the millions of 3D Gaussian points actually represent a limited number of unique semantics for a 3D scene, our LangSplatV2 assumes that each Gaussian can be represented as a sparse coding over a group of global basis vectors. Then, we derive that learning a high-dimensional feature field is equivalent to learning a 3D coefficient field and a global codebook. For each 3D Gaussian point, instead of augmenting a language feature, we
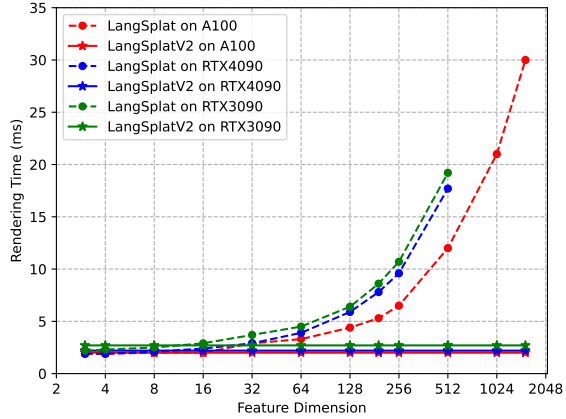


Figure 1: Feature rendering time comparison with different GPUs. Note that the less advanced GPUs (RTX 3090 and RTX 4090) cannot accommodate the LangSplat model with feature dimensions of 1024 or higher due to running out of memory.

learn sparse coefficients. During testing, we utilize the sparsity and propose an efficient sparse coefficient splatting method with CUDA optimization, which effectively decouples the rendering dimension with feature dimensions. As our method removes the autoencoder and directly learns the 3D field from high-dimensional 2D CLIP features, more accurate language fields are also obtained. Experimental results demonstrate that our LangSplatV2 not only achieves higher accuracy but also delivers a substantial speedup: 42 × faster than LangSplat for high-dimensional feature rendering and 47 × faster for open-vocabulary 3D querying.

## 2   Related Work

**3D Gaussian Splatting.** Recently, 3D Gaussian Splatting (3D-GS) [11] has received huge attention [16, 17] due to its significant speed advantage over neural radiance fields (NeRFs) [10]. As a fundamental 3D modeling technology, it's been widely used in many downstream tasks [18, 19, 20]. Yang *et al.* [21] extended 3D Gaussian Splatting to dynamic scenes by modeling deformable 3D Gaussians. Concurrently, 4D Gaussian Splatting [22] proposes to use multi-resolution voxel planes to model deformed Gaussians and attains high-quality video rendering results in real-time. Due to the explicit 3D modeling nature of 3D Gaussian Splatting, it has also been used for 3D editing. Gaussianeditor [23] offers swift and controllable 3D editing by tracing the editing target throughout the training process. Another concurrent work [24] edits 3D scenes using text instruction, which attains

two times faster training speed compared with Instruct-NeRF2NeRF [25]. Many methods [26, 27, 28] utilize 3D Gaussian Splatting for text-driven 3D generation. GSGEN [29] incorporates direct geometric priors from 3D Gaussians and obtains a text-to-3D generation model that captures high-frequency components. LangSplat [12] adopts 3D-GS to model a 3D language field and is $199 \times$ faster than the previous state-of-the-art method LERF. However, LangSplat is still not a real-time method, which undermines its broad application prospects.

**3D Language Field.** The concept of 3D language fields [30, 31, 32] has emerged as an intersection of vision-language models and 3D modeling, aiming to create interactive environments that can be manipulated and queried using natural language. LERF [13] explores 3D open-vocabulary text queries by embedding the CLIP feature into a 3D field. Liu *et al.* [33] also distilling knowledge from pre-trained foundation models like CLIP [8] and DINO [34] into NeRF [10]. It further proposes relevancy-distribution alignment and feature-distribution alignment losses to address the CLIP feature ambiguity issue. LangSplat improves LERF with 3D Gaussian Splatting and attains a significant speedup. There are also some other works [35, 36, 37] employing 3D Gaussian Splatting for 3D language field modeling, they usually suffer from imprecise language fields like LERF, which has been well addressed by LangSplat. Following LangSplat, GOI [38] proposes simplifying feature selection by dividing the feature space with a hyperplane, keeping only the features most relevant to the query. GAGS [39] enhances multiview consistency by linking SAM's prompt point density with camera distances and introducing an unsupervised granularity factor to selectively distill consistent 2D features. However, these methods both use a decoder to map the low-dimensional feature to high-dimensional space, which is a primary bottleneck of real-time query as shown in section 3.2.

## 3 Proposed Approach

### 3.1 Preliminaries

3D Gaussian Splatting employs point-based rendering technologies [40, 41] and models the scene geometry as a set of 3D Gaussian points. Each Gaussian point is associated with the following attributes: Gaussian center position $\mu \in \mathbb{R}^3$, a covariance matrix $\Sigma \in \mathbb{R}^{3 \times 3}$, color described $c \in \mathbb{R}^3$ by spherical harmonic (SH) coefficients, and opacity $o \in \mathbb{R}$. Each Gaussian point is represented as:

$$G(x) = \exp(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)). \tag{1}$$

The covariance matrix $\Sigma$ is further represented with a rotation matrix and a scaling factor matrix:

$$\Sigma = RSS^\top R^\top, \tag{2}$$

where $R \in \mathbb{R}^4$ denotes the rotation matrix and $S \in \mathbb{R}^3$ represents the scaling matrix.

Based on EWA volume splatting [42], 3D Gaussian Splatting projects the 3D Gaussian points onto a 2D image plane, which blends the ordered Gaussian points that overlap with the rendered pixel $v$:

$$C(v) = \sum_{i \in \mathcal{N}} c_i \alpha_i \prod_{j=1}^{i-1}(1 - \alpha_j), \tag{3}$$

where $\mathcal{N}$ represents the ordered Gaussian points within a tile, and $\alpha_i = o_i G_i^{2D}(v)$. The $G_i^{2D}(\cdot)$ means the projected 2D Gaussian distribution of the $i$-th 3D Gaussian point.

As 3D Gaussian Splatting exhibits excellent real-time rendering speed while maintaining high rendering quality even at 1080p resolution, many methods [22, 21, 43] have adopted it as a 3D scene modeling technology to accelerate the rendering speed. LangSplat aims to build a 3D language field to support open-vocabulary querying within 3D spaces. It extends each 3D Gaussian with a language embedding and supervises the 3D language Gaussians with 2D CLIP image features. LangSplat presents two main contributions to make it faster and more accurate. First, instead of using multiple patch-wise CLIP features with varying scales, which leads to imprecise and vague 3D language fields, LangSplat employs the CLIP features of SAM masks with three pre-defined SAM hierarchical semantic scales to obtain precise language fields with clear boundaries. Second, as CLIP embeddings are high-dimensional, directly splatting at CLIP feature space poses huge challenges in training memory and time cost. Therefore, LangSplat proposes a scene-specific autoencoder, which

Table 1: The stage-wise time cost (ms) for LangSplat and our improvements on the LERF dataset with one A100 GPU. LangSplat* is modified from LangSplat with simple engineering optimization. LangSplatV2 achieves a speed of up to 384.6 FPS for open-vocabulary 3D querying.

| Method | Rendering | Decoding | Post-Processing | Total | Speed (FPS) |
|--------|-----------|----------|-----------------|-------|-------------|
| LangSplat | 6.0 | 83.1 | 33.0 | 122.1 | 8.2 |
| LangSplat* | 2.0 | 83.1 | 0.5 | 85.6 | 11.7 |
| LangSplatV2 | **2.0** | **0.1** | **0.5** | **2.6** | **384.6** |

compresses the high-dimensional (512-D) CLIP features into low-dimensional (3-D) latent features. The rendering process of LangSplat follows:

$$\boldsymbol{F}(v) = \sum_{i \in \mathcal{N}} \boldsymbol{f}_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \tag{4}$$

where $\boldsymbol{f}_i \in \mathbb{R}^d$ represents the augmented $d$-dimensional latent feature at $i$-th 3D Gaussian point and $\boldsymbol{F}(v)$ denotes the rendered latent feature at pixel $v$. After obtaining the rendered feature, LangSplat uses the decoder $g_d(\boldsymbol{F}(v)) \in \mathbb{R}^D$ from the trained autoencoder $g$ to decode the latent feature back to CLIP feature space. In the end, the LangSplat attains an accurate 3D language field while being 199 $\times$ faster than the previous state-of-the-art method LERF [13].

## 3.2 Bottleneck Analysis

While LangSplat achieved significant speedup over other methods, it still cannot perform real-time open-vocabulary 3D querying at high resolution. However, there is a high demand for real-time 3D querying for many applications, such as Augmented Reality (AR) [44], intelligent robotics [13].

To further improve the query speed of LangSplat, we first analyze each step of the query process and identify the bottleneck. For a given text query, the inference of LangSplat can be divided into three stages: rendering, decoding, and post-processing. The rendering stage will render the learned 3D language Gaussians into a 2D image plane and get a rendered latent feature map $\boldsymbol{F} \in \mathbb{R}^{H \times W \times d}$, where $H, W$ denote the height and width of the 2D image size, respectively. As the rendered feature map only encodes language features in the low-dimensional latent space, a decoding stage is followed to obtain the feature map at CLIP space with a decoder $g_d(\boldsymbol{F}) \in \mathbb{R}^{H \times W \times D}$, where the decoder $g_d(\cdot)$ is implemented with an MLP. The last post-processing stage computes the relevancy score with the obtained $D$-dimensional feature following LERF [13] and remove some noise in the relevancy score map. Specifically, a mean filter is applied to the relevancy score map. Note that LangSplat adopts the three semantic levels introduced by SAM [45], so the above operations are repeated three times for three SAM semantic levels and the post-processing stage needs to select one semantic level for prediction with some specific strategies [12].

We test the stage-wise inference time of LangSplat on the LERF dataset with one A100 GPU. The results in Table 1 show that the rendering stage takes 6.0 ms, the decoding stage takes 83.1 ms, and the post-processing stage takes 33.0 ms. While all stages occupy a considerable amount of time, we can significantly improve them through two simple modifications. First, we notice that LangSplat implements the post-processing stage on the CPU, which could be slow when it involves mean filter operations. We implement the post-processing stage on the GPU and observe that the time cost of the post-processing stage is now negligible compared with other stages. The second modification is scale parallelization. LangSplat performs text querying on three semantic levels sequentially, meaning the rendering stage is performed three times separately. However, we can perform inference at three semantic scales in parallel. For example, instead of splatting $d$-dimensional features three times, we directly splat $3d$-dimensional features for once. We use LangSplat* to denote the version with these simple modifications. We observe that for LangSplat*, the rendering stage takes 2.0 ms, and the post-processing stage takes 0.5 ms. In these three stages, the decoding phase occupies 97.1% of the total query time, making it the primary bottleneck in achieving real-time 3D querying.

## 3.3 3D Sparse Coefficient Field

Different from 3D Gaussian Splatting, which renders 3-dimensional RGB color, 3D language fields need to splat high-dimensional features. Directly splatting high-dimensional features significantly
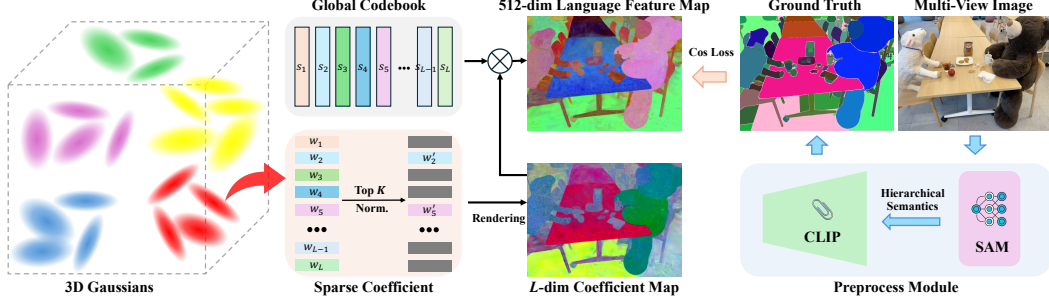
Figure 2: The framework of LangSplatV2. LangSplatV2 introduces a sparse coefficient for each Gaussian point and a shared global codebook for the entire scene.

decreases the rendering speed, as shown in Figure 1. To address this issue, existing methods usually model $d$-dimensional ($d << D$) latent features in 3D followed by either an online decoder [37] or an offline decoder [12] to decode latent features back to $D$-dimensional features in 2D image space. The high dimension gap between latent features and decoded features implies a heavyweight MLP is required to ensure the accuracy of the decoding stage, which becomes the primary speed bottleneck.

In LangSplatV2, instead of splatting the $d$-dimensional latent language features, we propose to model a 3D sparse coefficient field. As shown in Eq. 4, LangSplat assigns each Gaussian point with a language feature $\boldsymbol{f}_i$, which represents the semantics associated with the Gaussian point. As LangSplat could create millions of Gaussian points, there will be millions of unique language features. However, this is highly inefficient as the unique semantics within a scene are quite limited and much smaller than the number of Gaussian points. As a matter of fact, many Gaussian points share similar semantics. Therefore, we assume the language embedding of every Gaussian point within a scene can be represented as a sparse coding of $L$ global basis vectors $\mathcal{S} = [\boldsymbol{s}_1, \boldsymbol{s}_2, ..., \boldsymbol{s}_L]^\top \in \mathbb{R}^{L \times D}$. These $L$ basis vectors serve as the global codebook and each Gaussian point is computed by linearly combining a small number of local basis vectors. We assume that only $K$ ($K << L$) basis vectors from $L$ global codebook are used to represent the language embedding of a Gaussian point. Then we define $\boldsymbol{w}_i = [w_{i,1}, w_{i,2}, ..., w_{i,L}]^\top \in \mathbb{R}^{1 \times L}$ as the associated sparse coefficients for $i$-th Gaussian point, where $\sum_{l=1}^{L} w_{i,l} = 1$ and only $K$ elements are non-zero values while all other elements are zeros. With slightly abuse of the notion $\boldsymbol{f}_i$, the language embedding $\boldsymbol{f}_i \in \mathbb{R}^D$ of $i$-th Gaussian point is represented as:

$$\boldsymbol{f}_i = \boldsymbol{w}_i \mathcal{S} = \sum_{l=1}^{L} w_{i,l} \boldsymbol{s}_l. \tag{5}$$

If we directly render the $D$-dimensional language field without compressing CLIP features, we have:

$$\boldsymbol{S} = \sum_{i \in \mathcal{N}} \boldsymbol{w}_i \mathcal{S} \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \tag{6}$$

where $\boldsymbol{S}$ is the rendered $D$-dimensional CLIP feature. We set $e_i = \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j)$, then we have:

$$\boldsymbol{S} = \sum_{i \in \mathcal{N}} \boldsymbol{w}_i \mathcal{S} e_i = (\sum_{i \in \mathcal{N}} e_i \boldsymbol{w}_i) \mathcal{S}. \tag{7}$$

Eq. 7 shows that rendering the $D$-dimensional CLIP features is equivalent to first rendering the sparse coefficients $\boldsymbol{w}(i)$, and then performing a matrix multiplication with the global dictionary $\mathcal{S}$.

Therefore, we propose to learn a $L$-dimensional sparse coefficient for each Gaussian point and a global codebook with a size of $L$ basis vectors. To learn a sparse distribution for the coefficient $\boldsymbol{w}_i$, we first apply a softmax function to normalize the $L$-dimensional parameter, then preserve the top-$K$ values as non-zeros elements and set the remaining elements to zeros. We will re-normalize the top-$K$ values with their sum to ensure the sum of the coefficient equals one. The $L$-dimensional 3D sparse coefficient field and the $L$ basis vectors are jointly learned. Figure 2 visualizes the framework.

Compared to LangSplat, LangSplatV2 requires only a simple matrix multiplication after rendering the weight map, instead of relying on a heavyweight decoder, thus overcoming LangSplat's main inference

5

speed bottleneck. Furthermore, the $D$-dimensional global codebook eliminates reconstruction loss from dimensionality reduction, enabling better modeling of high-dimensional features in the CLIP space and improving query accuracy.

## 3.4 Efficient Sparse Coefficient Splatting

By learning a 3D sparse coefficient field, the MLP decoder is entirely removed, eliminating the associated computational overhead of the MLP. However, we still need to perform a $3L$-dimensional feature rendering and a matrix multiplication. Our experiments show that the time for matrix multiplication (in Table 1, listed as the decoding stage of LangSplatV2) is negligible compared to the rendering process. However, rendering high-dimensional features with dimensionality $L$ remains computationally demanding. To address this issue, we propose an efficient sparse coefficient splatting method with CUDA optimization. By exploiting the sparsity of the coefficient field, we can achieve $L$-dimensional feature rendering at the cost of only $K$-dimensional rendering, where $K \ll L$.

In the CUDA implementation of 3D Gaussian Splatting and LangSplat, each thread performs alpha-blending of the $|\mathcal{N}|$ ordered Gaussian points within a tile. With an $L$-dimensional rendering, each thread sequentially computes $L$ channels, leading to a computational complexity of $O(|\mathcal{N}|L)$. When $L$ becomes sufficiently large, this alpha-blending computation becomes the key bottleneck and significantly increases the overall computational overhead as $L$ grows, as shown in Figure 1.



To mitigate this, we utilize the sparse nature of the learned coefficient field during test-time querying, as shown in Figure 3. Although each Gaussian point has an $L$-dimensional coefficient, only $K$ dimensions contain non-zero values. Therefore we can only perform alpha-blending

Figure 3: Our efficient sparse coefficient splatting method accelerates the speed of alpha-blending by utilizing the property of the learned sparse coefficient field and neglecting zero elements.

for the non-zero elements, as the blending of zero elements does not alter the result. This effectively reduces the computational complexity from $O(|\mathcal{N}|L)$ to $O(|\mathcal{N}|K)$, with $K$ significantly smaller than $L$. In practice, we set $K = 4$, which allows us to simultaneously render three semantic scales with an effective feature rendering dimension in CUDA of only 12, yielding a high-quality feature map of $1,536$ dimensions. It makes the rendering highly efficient without compromising feature quality.

Specifically, each Gaussian's $L$-dimensional sparse coefficient $\boldsymbol{w}_i$ can be fully represented by its top-$K$ non-zero elements. We store these as two $K$-dimensional arrays for each Gaussian: top-$K$ indices and top-$K$ coefficients. The top-$K$ indices are the positions of the top-$K$ non-zero elements within the $L$-dimensional vector and the top-$K$ coefficients are the values of these top-$K$ non-zero elements. During rendering, each CUDA thread performs alpha-blending only for the $K$ non-zero coefficients. For each Gaussian point within the tile, the CUDA thread will access the indices and coefficients of the top-$K$ elements and perform weighted summation solely on these $K$-dimensional indices, avoiding computation on zero elements. In this way, LangSplatV2 can obtain high-dimensional ($D$) feature splatting results at the cost of only ultra-low-dimensional ($K$) feature splatting.

## 4 Experiments

### 4.1 Datasets and Details

**Datasets.** We evaluate our method on the LERF, 3D-OVS, and Mip-NeRF360 datasets. The LERF dataset [13], captured using the iPhone App Polycam, contains in-the-wild scenes. For the open-vocabulary 3D object localization task, we adopt the augmented localization annotations provided by LangSplat [12] on the LERF dataset. Additionally, we use the segmentation ground truth from LangSplat [12] for the open-vocabulary 3D segmentation task on LERF. Beyond LERF, we also conduct 3D segmentation experiments on the 3D-OVS and Mip-NeRF360 [46] datasets. The 3D-OVS dataset [33] includes a collection of long-tail objects captured in diverse poses and backgrounds.
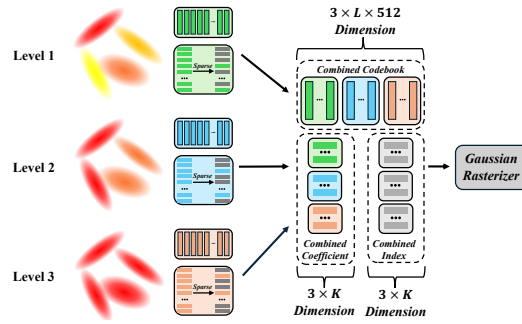
Table 2: Quantitative comparisons of open-vocabulary 3D object localization and 3D semantic segmentation on the LERF dataset. We report the mean accuracy (%) and the mean IoU scores (%).

| Method | 3D Object Localization | | | | | 3D Semantic Segmentation | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Ramen | Teatime | Kitchen | Figurines | Overall | Ramen | Teatime | Kitchen | Figurines | Overall |
| GS-Grouping [47] | 32.4 | 69.5 | 50.0 | 44.6 | 49.1 | 26.4 | 54.0 | 31.3 | 34.6 | 36.6 |
| LEGaussian [35] | 69.0 | 79.7 | 63.6 | 57.1 | 67.4 | 20.2 | 32.3 | 22.3 | 23.4 | 24.6 |
| GOI [38] | 56.3 | 67.8 | 68.2 | 44.6 | 59.2 | 33.7 | 55.8 | 54.5 | 23.9 | 42.0 |
| GAGS [39] | 69.0 | <u>88.1</u> | <u>90.9</u> | 78.6 | 81.7 | 46.8 | 60.3 | <u>55.8</u> | <u>53.6</u> | <u>54.1</u> |
| LangSplat [12] | <u>73.2</u> | <u>88.1</u> | **95.5** | <u>80.4</u> | **84.3** | <u>51.2</u> | <u>65.1</u> | 44.5 | 44.7 | 51.4 |
| LangSplatV2 | **74.7** | **93.2** | 86.4 | **82.1** | <u>84.1</u> | **51.8** | **72.2** | **59.1** | **56.4** | **59.9** |



Figure 4: Qualitative comparisons of open-vocabulary 3D object localization on the LERF dataset. The red points are the model predictions and the black dashed bounding boxes denote the annotations. We observe that LangSplatV2 generates better results than LangSplat.

Moreover, we evaluate our method on Mip-NeRF360, which consists of multi-view indoor and outdoor scene images, with segmentation labels annotated by GAGS [39]. For evaluation, we use localization accuracy for the 3D object localization task and report the average IoU scores for the open-vocabulary 3D segmentation task.

**Implementation Details.** Following LangSplat [12], we use the OpenCLIP ViT-B/16 model to extract CLIP features. We employ the ViT-H model for SAM [45] to segment images and obtain masks with three hierarchical semantics. The codebook size $L$ is set to 64 and the $K$ is set to 4. During test-time querying, we render three semantic scales simultaneously, leading to the actual rendering dimension of 12. The 3D Gaussians are first trained with RGB supervision for 30,000 iterations to reconstruct the RGB scene. Then we train another 10,000 iterations for the 3D sparse coefficient field by fixing all other 3D Gaussian parameters. All our experiments are conducted on one A100 GPU.

## 4.2 Quantitative Results

**Time Analysis.** To assess the speed improvement of LangSplatV2 over its LangSplat, we conducted a detailed time analysis using the LERF dataset, with the computations performed on a single A100 GPU. Table 1 presents a stage-wise breakdown of the time costs for LangSplat, LangSplat*, and LangSplatV2. LangSplat, our initial model, achieves a frame rate of 8.2 FPS. The LangSplat* improvements mainly focused on rendering and post-processing optimizations, which reduced the total time to 85.6 ms and increased the speed to 11.7 FPS. Here,

Table 3: Quantitative 3D semantic segmentation results on 3D-OVS, reported as mean IoU (%).

| Method | Bed | Bench | Room | Sofa | Lawn | Overall |
|---|---|---|---|---|---|---|
| FFD [1] | 56.6 | 6.1 | 25.1 | 3.7 | 42.9 | 26.9 |
| LERF [13] | 73.5 | 53.2 | 46.6 | 27.0 | 73.7 | 54.8 |
| 3D-OVS [33] | 89.5 | 89.3 | 92.8 | 74.0 | 88.2 | 86.8 |
| GS-Grouping [47] | 83.0 | 91.5 | 85.9 | 87.3 | 90.6 | 87.7 |
| LEGaussian [35] | 84.9 | 91.1 | 86.0 | 87.8 | 92.5 | 88.5 |
| GOI [38] | 89.4 | 92.8 | 91.3 | 85.6 | 94.1 | 90.6 |
| LangSplat [12] | <u>92.5</u> | <u>94.2</u> | <u>94.1</u> | <u>90.0</u> | <u>96.1</u> | <u>93.4</u> |
| LangSplatV2 | **93.0** | **94.9** | **96.1** | **92.3** | **96.6** | **94.6** |

Table 5: Ablation study on codebook size $L$. We report the average performance of the 3D object localization and 3D semantic segmentation tasks on the LERF dataset.

| $L$ | 32 | 64 | 128 |
|---|---|---|---|
| Accuracy (%) | 72.8 | 84.1 | 84.1 |
| IoU (%) | 53.9 | 59.9 | 60.5 |

Table 6: Ablation study on different $K$. We report the average performance of the 3D object localization and 3D semantic segmentation tasks on the LERF dataset.

| $K$ | 2 | 4 | 6 | 8 |
|---|---|---|---|---|
| Accuracy (%) | 79.4 | 84.1 | 84.4 | 83.6 |
| IoU (%) | 54.4 | 59.9 | 59.9 | 59.9 |

the rendering time was significantly cut to 2.0 ms, and the post-processing time was reduced to less than 1.0 ms with our proposed simple engineering modifications. In contrast, LangSplatV2 proposed an efficient sparse coefficient splatting method, that entirely removed the MLP decoder in the decoding stage. The only operation in the decoding stage is performing a matrix multiplication between the rendered $L$-dimensional coefficient and the global dictionary, which takes only 0.1 ms. While completely removing the MLP decoder, our LangSpaltV2 only needs 2.0 ms to render the sparse coefficient, leading to an overall 476.2 FPS for obtaining high-dimensional feature maps. In the end, LangSplatV2 achieved an exceptional reduction in total querying time to 2.6 ms. These results not only demonstrate a substantial enhancement in processing speed but also affirm the real-time capabilities of LangSplatV2 for high-resolution, complex 3D scene querying.

**Main Results on the LERF dataset.** Table 2 demonstrates the comparisons on the LERF dataset. For the open-vocabulary 3D object localization task, LangSplatV2 achieved the highest accuracy with notable improvements over previous methods for most scenes, including "ramen", "figurines", and "teatime". In the "Kitchen" scene, LangSplat achieved a higher accuracy than LangSplatV2. This difference is mainly due to the limited sample size (22 examples in totoal), with only a two-prediction gap (21 vs. 19). Given LangSplat's high variance in this scene, we ran it multiple times, obtaining an average of $17.5 \pm 1.8$ correct predictions, which suggests that LangSplatV2 provides more consistent performance advantage across different scenarios. Results on 3D semantic segmentation demonstrate that LangSplatV2 outperforms all other methods, and it consistently surpasses LangSplat across all scenes, highlighting the effectiveness of our proposed method.

**Main Results on the 3D-OVS dataset.** Table 3 shows the quantitative comparisons of open-vocabulary 3D semantic segmentation performance, on the 3D-OVS dataset across five test scenes. We see that LangSplatV2 achieves the highest overall mean IoU score of 94.6%. It significantly outperforms LangSplat across all different scenes, further validating its ability to build more accurate language fields.

Table 4: Quantitative 3D semantic segmentation results on Mip-NeRF360, reported as mean IoU (%).

| Method | Room | Counter | Garden | Bonsai | Overall |
|---|---|---|---|---|---|
| GS-Grouping [47] | 54.4 | 47.7 | 40.4 | 54.1 | 49.2 |
| LEGaussian [35] | 25.5 | 35.3 | 33.2 | 22.3 | 29.1 |
| GOI [38] | 60.3 | 46.6 | 59.8 | 67.3 | 58.5 |
| GAGS [39] | **65.2** | 61.1 | _61.2_ | _70.5_ | _64.5_ |
| LangSplat [12] | 53.2 | _68.8_ | 51.9 | 55.4 | 57.3 |
| LangSplatV2 | _64.3_ | **75.1** | **65.0** | **73.1** | **69.4** |

**Main Results on the Mip-NeRF360 dataset.** Table 4 presents the comparison results on the Mip-NeRF360 dataset. Notably, LangSplatV2 achieved a significant improvement in the overall average IoU score, reaching 69.4%, which is a marked advancement over LangSplat's 57.3%.

**Ablation Study.** In Table 5, we evaluate the effect of codebook sizes $L$ on the localization and segmentation tasks. We report the mean accuracy and IoU scores on the LERF dataset. Our results reveal that increasing L from 32 to 64 yields a notable improvement in both localization accuracy (72.8% to 84.1%) and segmentation IoU (53.9% to 59.9%). This suggests that a larger codebook size better captures the semantic fields in complex scenes, allowing for more precise language representations. However, further increasing $L$ results in only slightly increased performance, suggesting that $L = 64$ has already reached saturation on the LERF dataset. Table 6 further reports the effects of different $K$. We observe that $K = 4$ can saturate the performance of the Gaussian model, and a larger K will increase the rendering dimensions, thereby slowing down the rendering speed. More ablation studies and details can be found in the supplementary material.

**Discussion.** Recent work has also explored codebook-based language field representations, such as LEGaussians [35]. However, our approach departs significantly from theirs in several key aspects.
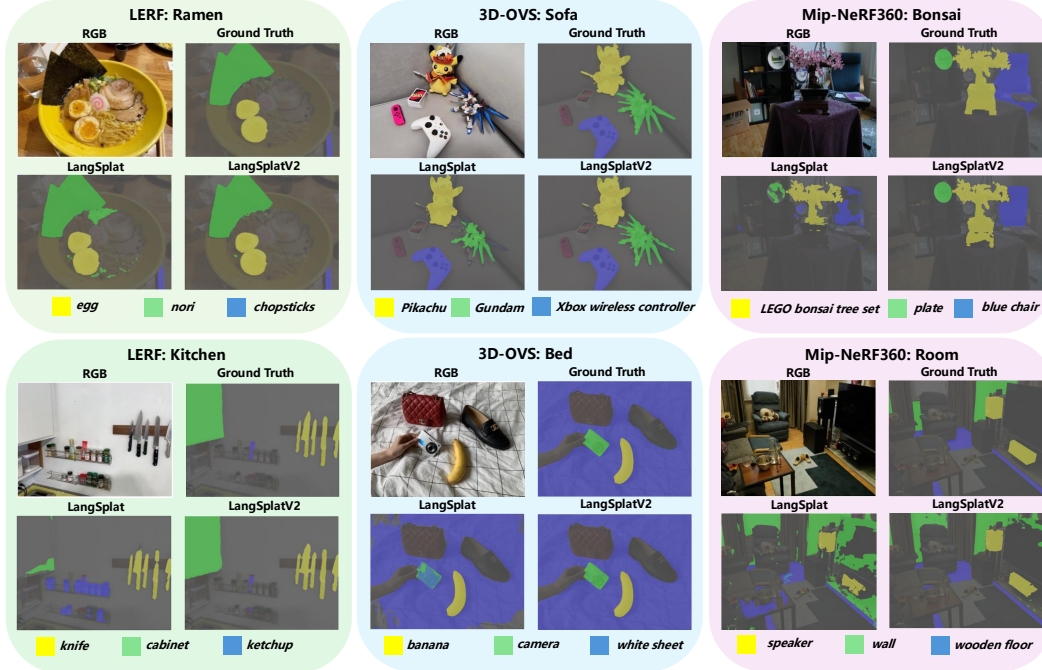
Figure 5: Qualitative comparisons of open-vocabulary 3D semantic segmentation on the LERF, Mip-NeRF360 and 3D-OVS dataset. We can see that our LangSplatV2 generates better masks than LangSplat, which shows the effectiveness of our LangSplatV2.

Table 7: Quantitative comparisons with LEGaussian on three benchmarks.

| Method | Segmentation IoU (%) | | | Query Time (ms) | | | GPU Memory (GB) | | |
|---|---|---|---|---|---|---|---|---|---|
| | LERF | 3DOVS | Mip. | LERF | 3DOVS | Mip. | LERF | 3DOVS | Mip. |
| LEGaussian [35] | 24.6 | 88.5 | 29.1 | 36.7 | 59.6 | 58.0 | 8.2 | **10.5** | **14.0** |
| LangSplatV2 | **59.9** | **94.6** | **69.4** | **2.6** | **4.8** | **3.8** | **7.2** | 11.6 | 16.9 |

First, LEGaussians construct a 2D feature codebook from 2D images while LangSplatV2 directly learns a global codebook in 3D space, which is more efficient. Second, LEGaussians still rely on an MLP to compress features and do not exploit sparsity. In contrast, our method completely removes the MLP and further enforces and leverages sparsity in both representation and rendering. As shown in Table 7, our method achieves superior semantic segmentation accuracy, faster inference speed, and comparable memory usage across three benchmark datasets, which demonstrates the effectiveness of our approach. For a more detailed comparison and discussion, please refer to the Appendix.

### 4.3 Qualitative Results

Figure 4 visualizes the open-vocabulary 3D object localization results. We observe that our LangSpaltV2 can give more accurate localization predictions compared with LangSplat. For example, our LangSplatV2 can give the correct prediction for "pumpkin" while LangSplat entirely fails in this query. Figure 5 visualizes the open-vocabulary 3D segmentation results in three datasets. We observe that our LangSplatV2 can generate more accurate masks compared with LangSplat. For example, in the "Kitchen" scene, LangSplat predicts a noisy mask for the "ketchup" query, while our LangSplatV2 generates more precise and clean masks.

## 5 Conclusion

**Conclusion.** In this paper, we have presented LangSplatV2 for high-dimensional language feature splatting in 3D open-vocabulary querying, designed to overcome the speed limitations of LangSplat. Our analysis identified the decoding stage as the primary bottleneck in LangSplat. By modeling

each Gaussian point as a sparse code over a global dictionary, LangSplatV2 removes the need for a heavyweight decoder. We further proposed an efficient sparse coefficient splatting method with CUDA optimization, allowing LangSplatV2 to achieve high-dimensional feature splatting results at the cost of splatting ultra-low-dimensional features.

**Limitations and broader impacts.** While LangSplatV2 demonstrates improved performance and faster inference compared to LangSplat, it comes with increased training cost due to the need to construct high-dimensional semantic fields during training. Detailed analysis of the training cost is provided in the appendix. Furthermore, as our method directly inherits the semantic representations from the CLIP model, it also inherits its inherent biases. Addressing such biases remains an open research challenge and may benefit from recent advances in fairness-aware versions of CLIP.

# References

[1] Sosuke Kobayashi, Eiichi Matsumoto, and Vincent Sitzmann. Decomposing nerf for editing via feature field distillation. *Advances in Neural Information Processing Systems*, 35:23311–23330, 2022.

[2] Vadim Tschernezki, Iro Laina, Diane Larlus, and Andrea Vedaldi. Neural feature fusion fields: 3d distillation of self-supervised 2d image representations. In *2022 International Conference on 3D Vision (3DV)*, pages 443–453. IEEE, 2022.

[3] Alan B Craig. Understanding augmented reality: Concepts and applications. 2013.

[4] Ronald T Azuma. A survey of augmented reality. *Presence: teleoperators & virtual environments*, 6(4):355–385, 1997.

[5] Songyou Peng, Kyle Genova, Chiyu Jiang, Andrea Tagliasacchi, Marc Pollefeys, Thomas Funkhouser, et al. Openscene: 3d scene understanding with open vocabularies. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 815–824, 2023.

[6] Chenguang Huang, Oier Mees, Andy Zeng, and Wolfram Burgard. Visual language maps for robot navigation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 10608–10615. IEEE, 2023.

[7] Paola Cascante-Bonilla, Hui Wu, Letao Wang, Rogerio S Feris, and Vicente Ordonez. Simvqa: Exploring simulated environments for visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5056–5066, 2022.

[8] 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.

[9] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024.

[10] 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.

[11] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023.

[12] Minghan Qin, Wanhua Li, Jiawei Zhou, Haoqian Wang, and Hanspeter Pfister. Langsplat: 3d language gaussian splatting. In *CVPR*, 2024.

[13] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lerf: Language embedded radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19729–19739, 2023.

[14] Raymond Koon Chuan Koh, Henry Been-Lirn Duh, and Jian Gu. An integrated design flow in user interface and interaction for enhancing mobile ar gaming experiences. In *2010 IEEE International Symposium on Mixed and Augmented Reality-Arts, Media, and Humanities*, pages 47–52. IEEE, 2010.

[15] Robert Godwin-Jones. Augmented reality and language learning: From annotated vocabulary to place-based mobile games. 2016.

[16] Zehao Yu, Anpei Chen, Binbin Huang, Torsten Sattler, and Andreas Geiger. Mip-splatting: Alias-free 3d gaussian splatting. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19447–19456, 2024.

[17] Yang Fu, Sifei Liu, Amey Kulkarni, Jan Kautz, Alexei A Efros, and Xiaolong Wang. Colmap-free 3d gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20796–20805, 2024.

[18] Heng Yu, Joel Julin, Zoltán Á Milacski, Koichiro Niinuma, and László A Jeni. Cogs: Controllable gaussian splatting. *arXiv preprint arXiv:2312.05664*, 2023.

[19] Xian Liu, Xiaohang Zhan, Jiaxiang Tang, Ying Shan, Gang Zeng, Dahua Lin, Xihui Liu, and Ziwei Liu. Humangaussian: Text-driven 3d human generation with gaussian splatting. *arXiv preprint arXiv:2311.17061*, 2023.

[20] Zhicheng Lu, Xiang Guo, Le Hui, Tianrui Chen, Min Yang, Xiao Tang, Feng Zhu, and Yuchao Dai. 3d geometry-aware deformable gaussian splatting for dynamic view synthesis. *arXiv preprint arXiv:2404.06270*, 2024.

[21] Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. Deformable 3d gaussians for high-fidelity monocular dynamic scene reconstruction. *arXiv preprint arXiv:2309.13101*, 2023.

[22] Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Xinggang Wang. 4d gaussian splatting for real-time dynamic scene rendering. In *CVPR*, 2024.

[23] Yiwen Chen, Zilong Chen, Chi Zhang, Feng Wang, Xiaofeng Yang, Yikai Wang, Zhongang Cai, Lei Yang, Huaping Liu, and Guosheng Lin. Gaussianeditor: Swift and controllable 3d editing with gaussian splatting. *arXiv preprint arXiv:2311.14521*, 2023.

[24] Jiemin Fang, Junjie Wang, Xiaopeng Zhang, Lingxi Xie, and Qi Tian. Gaussianeditor: Editing 3d gaussians delicately with text instructions. *arXiv preprint arXiv:2311.16037*, 2023.

[25] Ayaan Haque, Matthew Tancik, Alexei A Efros, Aleksander Holynski, and Angjoo Kanazawa. Instruct-nerf2nerf: Editing 3d scenes with instructions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19740–19750, 2023.

[26] Shijie Zhou, Zhiwen Fan, Dejia Xu, Haoran Chang, Pradyumna Chari, Tejas Bharadwaj, Suya You, Zhangyang Wang, and Achuta Kadambi. Dreamscene360: Unconstrained text-to-3d scene generation with panoramic gaussian splatting. *arXiv preprint arXiv:2404.06903*, 2024.

[27] Jiaxiang Tang, Zhaoxi Chen, Xiaokang Chen, Tengfei Wang, Gang Zeng, and Ziwei Liu. Lgm: Large multi-view gaussian model for high-resolution 3d content creation. *arXiv preprint arXiv:2402.05054*, 2024.

[28] Jiawei Ren, Liang Pan, Jiaxiang Tang, Chi Zhang, Ang Cao, Gang Zeng, and Ziwei Liu. Dreamgaussian4d: Generative 4d gaussian splatting. *arXiv preprint arXiv:2312.17142*, 2023.

[29] Zilong Chen, Feng Wang, and Huaping Liu. Text-to-3d using gaussian splatting. *arXiv preprint arXiv:2309.16585*, 2023.

[30] William Shen, Ge Yang, Alan Yu, Jansen Wong, Leslie Pack Kaelbling, and Phillip Isola. Distilled feature fields enable few-shot language-guided manipulation. *arXiv preprint arXiv:2308.07931*, 2023.

[31] Yawar Siddiqui, Lorenzo Porzi, Samuel Rota Buló, Norman Müller, Matthias Nießner, Angela Dai, and Peter Kontschieder. Panoptic lifting for 3d scene understanding with neural fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9043–9052, 2023.

[32] Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and Andrew J Davison. In-place scene labelling and understanding with implicit scene representation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15838–15847, 2021.

[33] 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*, 36:53433–53456, 2023.

[34] 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.

[35] 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.

[36] Xingxing Zuo, Pouya Samangouei, Yunwen Zhou, Yan Di, and Mingyang Li. Fmgs: Foundation model embedded 3d gaussian splatting for holistic 3d scene understanding. *arXiv preprint arXiv:2401.01970*, 2024.

[37] 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. *arXiv preprint arXiv:2312.03203*, 2023.

[38] Yansong Qu, Shaohui Dai, Xinyang Li, Jianghang Lin, Liujuan Cao, Shengchuan Zhang, and Rongrong Ji. Goi: Find 3d gaussians of interest with an optimizable open-vocabulary semantic-space hyperplane, 2024.

[39] Yuning Peng, Haiping Wang, Yuan Liu, Chenglu Wen, Zhen Dong, and Bisheng Yang. Gags: Granularity-aware 3d feature distillation for gaussian splatting. *arXiv preprint arXiv:2412.13654*, 2024.

[40] Hanspeter Pfister, Matthias Zwicker, Jeroen Van Baar, and Markus Gross. Surfels: Surface elements as rendering primitives. In *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, pages 335–342, 2000.

[41] Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. Surface splatting. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 371–378, 2001.

[42] Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. Ewa volume splatting. In *Proceedings Visualization, 2001. VIS'01.*, pages 29–538. IEEE, 2001.

[43] Zi-Xin Zou, Zhipeng Yu, Yuan-Chen Guo, Yangguang Li, Ding Liang, Yan-Pei Cao, and Song-Hai Zhang. Triplane meets gaussian splatting: Fast and generalizable single-view 3d reconstruction with transformers. *arXiv preprint arXiv:2312.09147*, 2023.

[44] Priyanka Jain, Nivedita Bhirud, Subhash Tatale, Abhishek Kale, Mayank Bhale, Aakanksha Hajare, and NK Jain. Real-time interactive ar for cognitive learning. In *AI, IoT, Big Data and Cloud Computing for Industry 4.0*, pages 219–239. Springer, 2023.

[45] 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. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026, 2023.

[46] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5460–5469, 2022.

[47] Mingqiao Ye, Martin Danelljan, Fisher Yu, and Lei Ke. Gaussian grouping: Segment and edit anything in 3d scenes. In *European Conference on Computer Vision*, pages 162–179. Springer, 2024.

[48] Zhiwen Fan, Kevin Wang, Kairun Wen, Zehao Zhu, Dejia Xu, and Zhangyang Wang. Lightgaussian: Unbounded 3d gaussian compression with 15x reduction and 200+ fps. *arXiv preprint arXiv:2311.17245*, 2023.

[49] Simon Niedermayr, Josef Stumpfegger, and Rüdiger Westermann. Compressed 3d gaussian splatting for accelerated novel view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10349–10358, 2024.

[50] Jin-Chuan Shi, Miao Wang, Hao-Bin Duan, and Shao-Hua Guan. Language embedded 3d gaussians for open-vocabulary scene understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5333–5343, 2024.

# A   Algorithm

The proposed efficient sparse coefficient splatting process is shown in Algorithm 1.

---
**Algorithm 1** Efficient Sparse Coefficient Splatting

---
    Initialize $\boldsymbol{W} \in \mathbb{R}^L$ to zero
    **for** Gaussian point $i \in \mathcal{N}$ **do**
        Retrieve top-$K$ indices $\{j_1, j_2, ..., j_K\}$ and corresponding coefficients $\{w_{i,j_1}, w_{i,j_2}, ..., w_{i,j_K}\}$
        **for** each $j \in \{j_1, ..., j_K\}$ **do**
            $\boldsymbol{W}[j] \mathrel{+}= w_{i,j} \cdot \alpha_i \cdot \prod_{m=1}^{i-1}(1 - \alpha_m)$
        **end for**
    **end for**

---

# B   Discussion

LangSplatV2 proposes to model the high-dimensional language features of each Gaussian point as a sparse code within a global dictionary. While similar concepts, such as quantizing Gaussian point attributes (*e.g.*, SH coefficients [48], opacity [49], and semantic features [50]) into a codebook, have been explored in the domains of compression [49] and language embedding [50], our method is fundamentally different.

In the compression domain, existing approaches [49, 48] first learn unique attributes for each Gaussian point, which are then quantized into a codebook using a quantizer. Extending this concept to our scenario would require learning 512-dimensional semantic features for each Gaussian point before quantization. This process significantly increases memory consumption and training time, often exceeding the capacity of GPUs like the NVIDIA 3090 or 4090, leading to out-of-memory (OOM) errors. Moreover, even with advanced GPUs that can handle the training, the rendering process would still involve managing 512-dimensional features. As shown in Figure 1 of the main paper, this dramatically reduces rendering speed compared to LangSplatV2.

In the context of language embedding, LEGaussians [50] adopt a codebook-based approach to accelerate training. However, our method is significantly different from LEGaussians in three key aspects: 1) Global 3D Codebook vs. 2D Codebook: LEGaussians first train a 2D codebook by quantizing features in 2D images, then learn a 3D model to predict the one-hot class category indicating the index of the codebook. In contrast, LangSplatV2 learns a global 3D codebook shared among all 3D Gaussian points. By avoiding the additional reconstruction errors caused by first quantizing features in the 2D image plane, our approach reduces overall reconstruction error and more efficiently utilizes the codebook in 3D space. 2) Eliminating the MLP Bottleneck: LEGaussians rely on an MLP to project each Gaussian point's features into a one-hot class category representing the index of the codebook. Our method removes the need for an MLP entirely, thereby eliminating the speed bottleneck and improving computational efficiency. 3) Efficient Sparse Coefficient Splatting: LangSplatV2 introduces Efficient Sparse Coefficient Splatting, which reduces the effective rendering dimensions to $K$. In contrast, LEGaussians renders with the full feature dimensions of Gaussian points, which inherently slows down the rendering process.

Furthermore, due to the reliance on an MLP decoder, LEGaussians still operate in a lower-dimensional feature space, which limits the expressiveness of its learned representations. Additionally, its two-stage 2D codebook approach fails to effectively incorporate 3D priors, resulting in a suboptimal learned field. In contrast, LangSplatV2 directly models high-dimensional features within a global 3D codebook, effectively capturing the underlying 3D structures. As evidenced by the results in Table 7 of the main paper, our approach outperforms LEGaussians across multiple benchmarks while being significantly faster, demonstrating superior scene reconstruction quality and more accurate language-grounded 3D representations.

To conclude, LangSplatV2 demonstrates a significant leap forward in the efficiency and scalability of modeling high-dimensional language features in 3D scenes. By leveraging the proposed 3D sparse coefficient fields and efficient sparse coefficient splatting techniques, our method reduces computational overhead while maintaining high fidelity in 3D language scene rendering.

# C  More Ablation Study

**Time Analysis.** We conduct ablation studies on the LERF [13] dataset and report the render and decode speed (ms per query) in Table 8. We test the speed on one A100 GPU. The baseline is LangSplat, which renders three-level 3-dimensional language features separately. Rendering three semantic levels in parallel can reduce the rendering time from 6.0 ms to 2.0 ms per query. The sparse coefficient field can significantly speed up the decoding stage from 83.1 ms/q to 0.1 ms/q by changing the heavy weight decoder to a simple matrix multiplication operation. Efficient sparse coefficient splatting method effectively transformed the 192-dimensional rendering into 12-dimensional rendering, reducing the rendering time from 5.3 ms to 2.0 ms per query through CUDA optimization.

Table 8: Speed ablation study on the LERF dataset. Parallel means rendering three semantic levels in parallel, Sparse means 3D sparse coefficient field, Efficient means efficient sparse coefficient splatting with CUDA optimization.

| Component | | | Speed (ms/q) | | |
|:---:|:---:|:---:|:---:|:---:|:---:|
| Parallel | Sparse | Efficient | render | decode | total |
| | | | 6.0 | 83.1 | 89.1 |
| ✓ | | | 2.0 | 83.1 | 85.1 |
| ✓ | ✓ | | 5.3 | 0.1 | 5.4 |
| ✓ | ✓ | ✓ | **2.0** | **0.1** | **2.1** |

$L$ **and** $K$**.** In Table 9 and Table 10, we show the ablation study results on each scene in the LERF [13] dataset. As we can see, for some complex scenarios, such as Kitchen and Figurines, larger $L$ and $K$ can lead to better performance. However, for all scenarios, the model's performance is already good enough with the settings of $L = 64$ and $K = 4$, and increasing $L$ and $K$ will lead to increased computational resources and time cost for training and inference. Therefore, we set $L = 64$ and $K = 4$ to balance performance and efficiency.

Table 9: Ablation study on codebbok size $L$. We report the localization and the segmentation performance on the LERF dataset.

| | $L$ | 32 | 64 | 128 |
|:---|:---|:---:|:---:|:---:|
| Ramen | Accuracy (%) | 70.4 | **74.7** | <u>74.7</u> |
| | IoU (%) | 50.5 | **51.8** | <u>51.4</u> |
| Teatime | Accuracy (%) | 84.8 | **93.2** | <u>91.5</u> |
| | IoU (%) | 65.7 | **72.2** | <u>69.8</u> |
| Kitchen | Accuracy (%) | 68.2 | <u>86.4</u> | **86.4** |
| | IoU (%) | 54.2 | <u>59.1</u> | **63.2** |
| Figurines | Accuracy (%) | 67.9 | <u>82.1</u> | **83.9** |
| | IoU (%) | 45.3 | <u>56.4</u> | **57.6** |
| Overall | Accuracy (%) | 72.8 | **84.1** | <u>84.1</u> |
| | IoU (%) | 53.9 | <u>59.9</u> | **60.5** |

**Training Cost.** Table 11 shows the comparison of training cost on the LERF [13] dataset with one A100 GPU. Compared with LEGaussians [35] and LangSplat [12], LangSplatV2 comes with increased training cost due to the need to construct high-dimensional semantic fields during training. Although our LangSplatV2 incurs higher training costs compared to LangSplat and LEGaussians, our primary focus in this paper is on improving the model's test-time performance for deployment. As demonstrated in the paper, LangSplatV2 achieves significantly higher inference speed—47× faster than LangSplat and 14× faster than LEGaussians on the LERF dataset—under comparable memory consumption, while also yielding higher segmentation accuracy. Furthermore, compared with naive LangSplat (directly trained with 512 dimensions), LangSplatV2 is more efficient with the proposed global codebook and sparse coefficient field.

Table 10: Ablation study on different $K$. We report the localization and the segmentation performance on the LERF dataset.

| | $K$ | 2 | 4 | 6 | 8 |
|---|---|---|---|---|---|
| Ramen | Accuracy (%) | 71.8 | **74.7** | 71.8 | <u>73.2</u> |
| | IoU (%) | 50.2 | **51.8** | <u>51.4</u> | 51.1 |
| Teatime | Accuracy (%) | 91.5 | **93.2** | <u>91.5</u> | 91.5 |
| | IoU (%) | 68.8 | **72.2** | <u>70.7</u> | 70.6 |
| Kitchen | Accuracy (%) | 77.3 | 86.4 | **95.5** | <u>90.9</u> |
| | IoU (%) | 48.9 | 59.1 | <u>59.8</u> | **60.1** |
| Figurines | Accuracy (%) | 76.8 | **82.1** | <u>78.6</u> | 78.6 |
| | IoU (%) | 49.6 | 56.4 | **57.7** | <u>57.6</u> |
| Overall | Accuracy (%) | 79.4 | <u>84.1</u> | **84.4** | 83.6 |
| | IoU (%) | 54.4 | **59.9** | <u>59.9</u> | 59.9 |

Table 11: Quantitative comparisons of training cost on the LERF dataset.

| Method | Training Time (h) | Training GPU Memory (GB) |
|---|---|---|
| LangSplat [12] | 1.0 | 6.2 |
| LEGaussian [35] | 1.3 | 11 |
| LangSplatV2 | 3.0 | 21.2 |
| LangSplat [12](512 Dim) | 45.0 | 47.4 |

## D  More Visualization Results

**Open-vocabulary 3D Object Localization.** We visualize more examples on the LERF [13] and Mip-NeRF360 [13] datasets for open-vocabulary 3D object localization in Fig 6.

**Open-vocabulary 3D Semantic Segmentation.** We visualize more examples on the LERF [13], Mip-NeRF360 [46] and 3D-OVS [33] datasets for open-vocabulary 3D semantic segmentation in Fig 7.
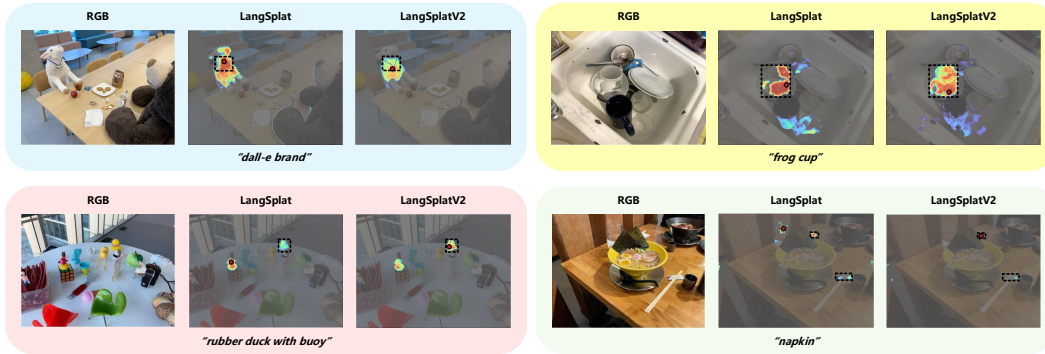
Figure 6: More qualitative comparisons of open-vocabulary 3D object localization on the LERF dataset. The red points are the model predictions and the black dashed bounding boxes denote the annotations. We observe that LangSplatV2 generates better results than LangSplat.
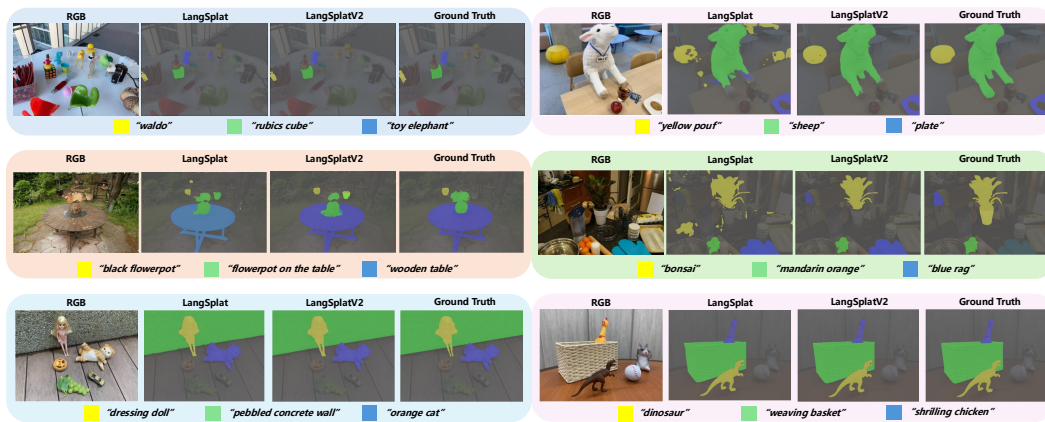


Figure 7: More qualitative comparisons of open-vocabulary 3D semantic segmentation on the LERF, Mip-NeRF360 and 3D-OVS dataset. We can see that our LangSplatV2 generates better masks than LangSplat, which shows the effectiveness of our LangSplatV2.