# FAST-Splat: Fast, Ambiguity-Free Semantics Transfer in Gaussian Splatting

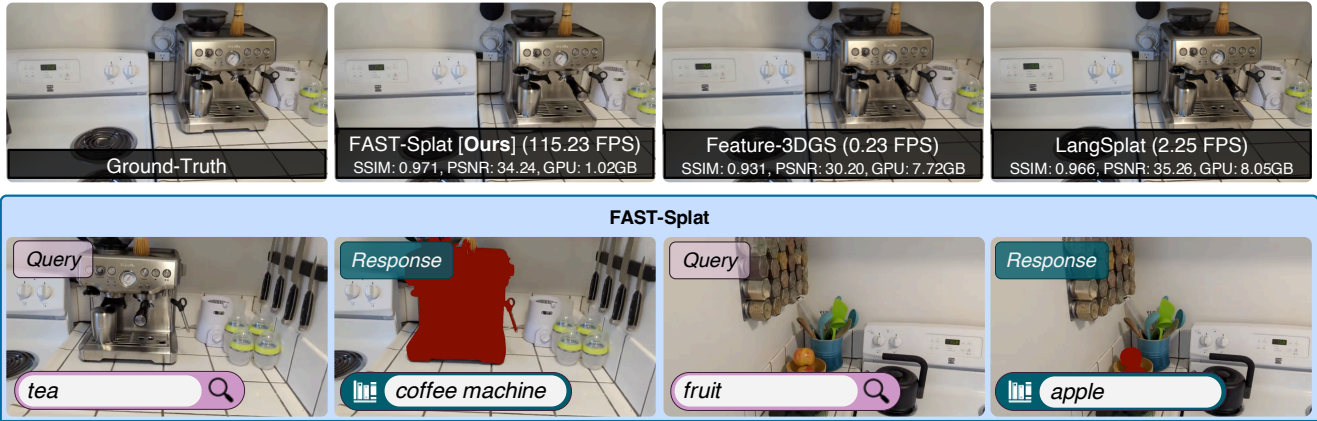Ola Shorinwa        Jiankai Sun        Mac Schwager
Stanford University

Figure 1. FAST-Splat enables fast, ambiguity-free semantic Gaussian Splatting, achieving 4x to 6x faster training times, 18x to 75x faster rendering speeds, and 3x lower GPU memory usage compared to prior work. FAST-Splat resolves language/scene-attributed ambiguity in object localization, providing the precise semantic label of objects when given a user query, e.g., in the illustrated cases with a related but different object, FAST-Splat informs a user that it found a *coffee machine* given the query: "tea," and an *apple* given the query: "fruit."

## Abstract

*We present FAST-Splat for fast, ambiguity-free semantic Gaussian Splatting, which seeks to address the main limitations of existing semantic Gaussian Splatting methods, namely: slow training and rendering speeds; high memory usage; and ambiguous semantic object localization. In deriving FAST-Splat, we formulate open-vocabulary semantic Gaussian Splatting as the problem of extending closed-set semantic distillation to the open-set (open-vocabulary) setting, enabling FAST-Splat to provide precise semantic object localization results, even when prompted with ambiguous user-provided natural-language queries. Further, by exploiting the explicit form of the Gaussian Splatting scene representation to the fullest extent, FAST-Splat retains the remarkable training and rendering speeds of Gaussian Splatting. Specifically, while existing semantic Gaussian Splatting methods distill semantics into a separate neural field or utilize neural models for dimensionality reduction, FAST-Splat directly augments each Gaussian with specific semantic codes, preserving the training, rendering, and memory-usage advantages of Gaussian Splatting over neural field methods. These Gaussian-specific semantic codes, together with a hash-table, enable semantic similarity to be measured with open-vocabulary user prompts and further enable FAST-Splat to respond with unambiguous semantic object labels and 3D masks, unlike prior methods. In experiments, we demonstrate that FAST-Splat is 4x to 6x faster to train with a 13x faster data pre-processing step, achieves between 18x to 75x faster rendering speeds, and requires about 3x smaller GPU memory, compared to the best-competing semantic Gaussian Splatting methods. Further, FAST-Splat achieves relatively similar or better semantic segmentation performance compared to existing methods. After the review period, we will provide links to the project website and the codebase.*

{shorinwa,jksun,schwager}@stanford.edu

## 1. Introduction

Research breakthroughs in vision-language foundation models have underpinned the remarkable performance of state-of-the-art object detection and classification [24, 50], segmentation [15, 29], and image captioning [27, 43]. In general, these models learn useful multi-modal image-language representations, entirely supervised with 2D image-text pairs, all within a shared representation space. Recent work has shown that grounding the semantic knowledge encoded by

vision-language foundation models in 3D can be useful in improving semantic segmentation and localization [12, 35] (especially in models that take in low-resolution images as input), enabling applications such as 3D scene-editing [14, 39] and robotic manipulation [28, 36]. Although these methods demonstrate compelling object segmentation performance from open-vocabulary queries, these methods fail to disambiguate between the semantic classes of unique objects that are semantically-similar to the natural language query [12]. For example, when a user queries for "tea," existing semantic Gaussian Splatting methods localize the coffee machine without providing any clarifying information on the identity of the localized object, suggesting that "tea" exists in the scene, when the contrary is true. This ambiguity generally arises from the fact that many vision-language foundation methods, e.g., CLIP, utilize a bag-of-words approach in measuring the semantic similarity between image-language queries. Moreover, many semantic methods for Gaussian Splatting are notably slow, requiring a significant amount of memory and computation time for training and inference. We seek to derive a semantic Gaussian Splatting method that addresses these limitations.

In this work, we propose FAST-Splat, a *Fast, Ambiguity-free Semantics Transfer* method for Gaussian Splatting. FAST-Splat generates fine-grained semantic localization, with highly specific semantic classes even in the presence of ambiguous natural-language prompts, and as its name implies, runs much faster than existing semantic Gaussian Splatting methods. FAST-Splat utilizes the vision-language foundation model CLIP [27] to extend closed-set semantic Gaussian Splatting to an open-world setting, while retaining the non-ambiguity of localized objects inherent in closed-set methods. Unlike existing semantic Gaussian Splatting methods which rely on neural-field models for semantic distillation, FAST-Splat directly augments each Gaussian with specific semantic codes, without any neural model. This key insight not only enables FAST-Splat to exploit the explicit scene-representation of Gaussian Splatting for faster training and rendering speeds, but also enables FAST-Splat to provide unambiguous semantic object labels in response to open-vocabulary user queries, e.g., in semantic object localization, upon identifying semantically-similar objects using the semantic codes associated with each Gaussian, along with the semantic hash-table.

We compare FAST-Splat to existing semantic Gaussian Splatting methods [26, 49]. FAST-Splat achieves 4x to 6x faster training times and 18x to 75x faster rendering speeds, while requiring 3x smaller GPU memory. Meanwhile, FAST-Splat achieves competitive semantic object localization performance, outperforming prior work in some scenes, while achieving the second-best performance in the other scenes. Notably, FAST-Splat enables semantic disambiguation in object localization, even with ambiguous user-provided natural-language queries, which we illustrate in Figure 1. In contrast to prior work, FAST-Splat provides a clarifying semantic label to each localized object, disambiguating the semantic identity of the localized object. For example, in Figure 1, when prompted with the query "tea," FAST-Splat notifies a user that it localized a *coffee machine*, and an *apple* when given the query: "fruit." Such disambiguation can be critical to enabling interesting downstream applications, e.g., in robotics.

To summarize our contributions:

- We introduce FAST-Splat, a Fast, Ambiguity-Free Semantic Gaussian Splatting method that enables notably faster (4x to 6x faster) language-semantics grounding in 3D scenes, distilled from 2D vision models.
- FAST-Splat resolves semantic ambiguity in semantic object localization, arising from ambiguous user-provided natural-language queries or inherent scene ambiguity, enabling the identification of the precise semantic label of relevant objects.
- FAST-Splat achieves superior rendering speeds for semantic Gaussian Splatting, about 18x to 75x faster rendering speeds compared to prior work.

## 2. Related Work

We provide a detailed review of existing work on open-vocabulary object detection and segmentation, along with follow-on work on grounding 2D object detection and segmentation in 3D by leveraging high-fidelity 3D scene representations.

**Open-Vocabulary Object Detection and Segmentation.** The introduction of the vision transformer (ViT) architecture [6] has spurred rapid research advances in many classical computer vision tasks, such as object detection and segmentation. However, like their convolutional neural network (CNN) counterparts [8, 31], early ViT models [2, 45] were primarily trained to detect objects within a fixed number of classes, defined as closed-set object detection. Vision-image foundation models, e.g., CLIP [27] and transformer-based language encoders, such as BERT [5], have enabled the detection of objects through open-vocabulary natural-language prompts, without relying on a fixed number of object classes, referred to as open-set object detection. Open-vocabulary object detection methods [7, 17, 22, 24, 44] fuse the image and language embeddings from vision or language foundation models using a decoder transformer architecture with self-attention to detect objects in a query image. GLIP [17] and OWL-ViT [24] train an image and text encoder jointly to learn useful representations, which are fed to a decoder for 2D grounding. While GLIP and OWL-ViT require a pretraining phase on large-scale internet image-text pairs, GroundingDINO [22] leverages pre-trained text encoder BERT [5], but otherwise utilizes a similar model architecture.

Zero-shot object segmentation methods [16, 21, 29, 41] leverage an encoder-decoder model architecture to generate a semantic mask of all instances of objects in an image. Mask DINO [16] introduces a mask prediction branch to DINO [45], extending it to object segmentation. Like GLIP, ViL-Seg [21] trains an image and text encoder together with a clustering head on image-caption pairs gathered on the internet to segment objects, enabling segmentation from natural-language prompts. SAM-2 [29] utilizes an image encoder with memory attention and a memory bank to enable video segmentation (in addition to image segmentation). However, SAM requires a mask, bounding box, or pixel-coordinate prompt for segmentation. GroundedSAM [32] extends SAM-2 to open-vocabulary queries by computing the bounding box associated with a given natural-language query, which is passed as an input prompt to SAM-2. Other methods [15, 18, 46] leverage pre-trained image and text encoders, which are fine-tuned in some cases, to enable open-vocabulary object segmentation.

**Radiance Fields.** Neural Radiance Fields (NeRFs) [1, 23, 25] represent a scene as a color and density field, parameterized by multi-layer perceptrons. NeRFs generate a photorealistic reconstruction of a scene, capturing high-fidelity visual details. However, high-quality NeRFs typically require long training times and achieve slow rendering rates. Gaussian Splatting [11] addresses these limitations by representing a scene using ellipsoidal primitives with visual attributes, such as opacity and color using spherical harmonics, in addition to the spatial and geometric attributes, such as the ellipsoid's mean and covariance. Radiance fields can be trained entirely from RGB images, making them suitable for many vision tasks.

**Grounding Language in 3D.** To bridge the gap between 2D object detection and 3D object localization, prior work has examined grounding semantic embeddings from pre-trained image and text encoders in the 3D world. CLIP-Fields [35] learns an implicit spatial mapping from 3D points to CLIP image embeddings and Sentence-Bert [30] text embeddings using back-projected RGB-D points. VLMaps [10] and NLMap [3] distill vision-language features into grid-based maps. The sparsity of these maps limits the resolution of the open-vocabulary segmentation results achievable by these methods. For high-resolution object segmentation, Distilled Features Fields (DFF) [14], CLIP-NeRF [39], and LERF [12] ground CLIP [27], LSeg [15], and DINO [45] features in a NeRF representation of the scene. Moreover, similar techniques have been employed in [9, 26, 49, 51] to distill semantic features into Gaussian Splatting representations. The resulting distilled radiance fields (DRFs) enable dense open-vocabulary object segmentation in 3D and have been applied to scene-editing [14, 39] and robotic manipulation [28, 36, 37, 48]. Despite their success in object segmentation, existing DRFs fail to precisely identify the semantic

class of the segmented object, e.g., for a given query, these methods localize all objects related to the query without resolving the unique semantic identity of each object. This ambiguity creates a gap in scene understanding, which is often required in many downstream tasks, e.g., in robotic exploration or manipulation. In this work, we introduce a DRF with fine-grained semantics, enabling precise semantic identification of objects in a scene, disambiguating language-guided semantic object localization, and addressing existing challenges.

## 3. Preliminaries

Here, we provide some background information on Gaussian Splatting. Gaussian Splatting [11] represents a scene using ellipsoidal primitives, with ellipsoid $i$ parameterized by its mean $\mu_i \in \mathbb{R}^3$; covariance $\Sigma = RSS^\mathsf{T}R^\mathsf{T} \in \mathbb{R}^{3 \times 3}$, where $R$ denotes a rotation matrix, and $S$ denotes a diagonal scaling matrix; opacity $\alpha$; and spherical harmonics (SH). Gaussian Splatting utilizes a tile-based rasterization procedure for efficient real-time rendering rates, outperforming NeRF-based representations, with the color $C$ of each pixel given by the $\alpha$-based blending procedure:

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

where $c_i$ denotes the color of each ellipsoid and $\alpha_i$ denotes the per-ellipsoid opacity multiplied by the probability density of the 2D Gaussian associated with the ellipsoid. Upon initialization from a sparse point cloud, which can be obtained from structure-from-motion, the attributes of each ellipsoid are optimized through gradient descent on a dataset of RGB images. Distilled feature fields build upon NeRFs and Gaussian Splatting, leveraging the same volumetric rendering and tile-based rasterization for NeRFs and Gaussian Splatting, respectively, to generate 2D feature maps from the 3D feature fields.

## 4. Fast, Ambiguity-Free Semantic Transfer

Now, we discuss FAST-Splat, a fast, fine-grained semantic distillation method for grounding 2D language semantics in 3D, enabling precise identification of the semantic label of objects in the 3D scenes. We build off 3D Gaussian Splatting to obtain a high-fidelity, photorealistic scene reconstruction entirely from RGB images. Given a dataset of images $\mathcal{D}$, Gaussian Splatting optimizes the mean, covariance, opacity, and SH parameters of each ellipsoid comprising the representation via the loss function:

$$\mathcal{L}_{\mathrm{gs}} = (1 - \lambda) \sum_{\mathcal{I} \in \mathcal{D}} \left\| \mathcal{I} - \hat{\mathcal{I}} \right\|_1 + \lambda \mathcal{L}_{\mathrm{D-SSIM}}, \qquad (2)$$

where $\lambda \in [0, 1]$ denotes a weighting parameter between the $\ell_1$-loss term and the differentiable structural similarity index

measure (D-SSIM) loss term, and $\mathcal{I}$ and $\hat{\mathcal{I}}$ denote the ground-truth and predicted RGB image, respectively. The Gaussian Splatting model is trained using the Adam optimizer [13], a gradient-based optimization method.

Almost all existing semantic distillation methods rely on neural architectures for grounding language in 3D in NeRFs [12, 14, 39] and Gaussian Splatting scenes [9, 26, 49, 51]. Although effective in enabling open-set semantic object localization, neural-based approaches introduce significant computational bottlenecks into the training pipeline of NeRF and Gaussian Splatting scenes. Moreover, when utilized in Gaussian Splatting, neural-based semantic distillation essentially erases the computational gains, e.g., faster training times and rendering rates, afforded by Gaussian Splatting when compared to NeRFs.

Via FAST-Splat, we seek to derive a method for open-set semantic distillation, while preserving the faster training times and rendering rates achieved by the underlying Gaussian Splatting method. To achieve this goal, we revisit the fundamental insight exploited by the authors of the original Gaussian Splatting work [11], namely: that only representing occupied space explicitly can lead to significant speedup in training and rendering times without any notable degradation in the visual fidelity of the representation. Specifically, Gaussian Splatting harnesses the explicit form and sparsity of an ellipsoidal-primitive-based scene representation to achieve faster training times and rendering times compared to NeRFs, the state-of-the-art at the time, which relied on implicit fields for the color and density of the scene. By learning neural encoders for dimensionality reduction or neural fields for language semantics, existing semantic distillation methods for Gaussian Splatting fail to take advantage of this insight. In contrast, in FAST-Splat, we leverage this insight to speed up training and rendering times, by eschewing neural models. In particular, we augment the attributes of each ellipsoid with a semantic parameter, associated with its semantic identity $\beta \in \mathbb{R}^3$, which is trained simultaneously with the ellipsoid's other attributes.

**Closed-Set Semantic Gaussian Splatting.** We formulate the open-set semantic distillation problem as the problem of extending closed-set semantic distillation (segmentation) methods to the open-set regime. To obtain the closed-set of semantic categories, we feed each image $I \in \mathcal{D}$ into an object detector, e.g., YOLO [40], DETR [2], and Efficient DETR [42], to generate a list of objects present in the image. Although closed-set object detectors generally exhibit remarkable performance, these detectors occasionally fail to identify all instances of an object in an image, given that the size of their dictionary of detectable objects is often limited. To alleviate this issue, we further process the outputs of the object detector using an open-set object detector, e.g., GroundingDINO [22], with the input prompt given by the classes of objects identified by the object detector. In prelim-

inary results, we find that this additional data pre-processing stage improves the robustness of the data pre-processing pipeline, by eliminating false detections made by the closed-set object detector.

At this stage, each image has a list of constituent objects and bounding boxes specifying the location of these objects in the image. However, high-fidelity semantic distillation generally requires pixel-wise semantic classes. To generate dense semantic maps, we utilize SAM-2 [29], passing in each object class and its associated bounding box. We maintain a dictionary (hash table), representing the closed-set of objects identified in the image. Now, we can distill the pixel-wise semantic information into the underlying Gaussian Splatting model. To predict the semantic class of each ellipsoid in the scene, we apply an affine transformation $\mathcal{W}$ to the semantic attribute $\beta$ of each ellipsoid, mapping $\beta$ to an $N$-dimensional space, where $N$ denotes the size of the dictionary. Subsequently, we apply the softmax function to the resulting outputs to define a probability over the entries in the dictionary. We extend this procedure to view-based rendering. Given a camera pose for rendering, we utilize differentiable tile-based rasterization [11] to render a semantic feature map corresponding to the camera pose. We utilize the aforementioned procedure to transform the semantic feature of each pixel to a probability distribution over the entries in the closed-set of object classes. We optimize the semantic attributes using the multi-class cross-entropy loss function, given by:

$$\mathcal{L}_{\mathrm{ce}} = -\sum_{\mathcal{I} \in \mathcal{D}} \frac{1}{|\mathcal{T}|} \sum_{p \in \mathcal{T}} \gamma_p \log \left( \frac{\exp(\mathcal{F}_{\mathcal{I},p,c})}{\sum_{j=1}^{N} \exp(\mathcal{F}_{\mathcal{I},p,j})} \right), \quad (3)$$

where $\mathcal{T}$ represents the set of all indices of pixels in the rendered feature map $\mathcal{F}_{\mathcal{I}}$ associated with image $\mathcal{I}$ (after applying the affine transformation $\mathcal{W}$), $c$ denotes the true class of the pixel, and $\gamma_p$ represents the weight associated with pixel $p$. We augment the set of object classes for each view with an *undetected* class label and assign a lower weight to pixels with this class label, i.e., pixels that were not segmented during the data pre-processing phase. This augmentation step serves as a regularization term for the stability of undetected pixels. We train all attributes of the Gaussian Splatting representation, e.g., the geometric, visual, and semantic attributes and the parameters in $\mathcal{W}$, simultaneously using the loss function: $\mathcal{L}_{\mathrm{sgs}} = \mathcal{L}_{\mathrm{gs}} + \mathcal{L}_{\mathrm{ce}}$, with $\mathcal{L}_{\mathrm{gs}}$ defined in (2).

**Open-Vocabulary Semantic Gaussian Splatting.** The resulting distillation method enables object localization from a pre-defined set of object classes. However, it does not support open-ended queries from users, which represents an important practical use-case. Here, we extend the capability of the distillation method from the closed-set setting to the open-vocabulary setting, leveraging open-source
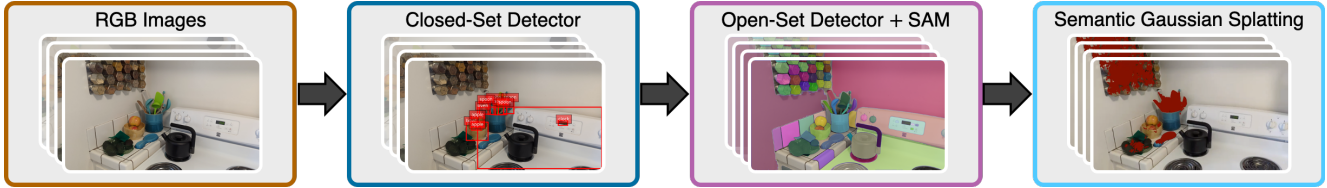
Figure 2. Unlike prior semantic Gaussian Splatting methods, FAST-Splat trains a non-neural-based semantic Gaussian Splatting model to achieve much faster training and rendering times by extending closed-set semantic segmentation to the open-vocabulary setting. FAST-Splat distills language semantics from a closed-set object detector, e.g., YOLO, augmented with an open-set object detector, e.g., GroundingDINO, and a image segmentation model, e.g., SAM-2.

pre-trained text encoders, e.g., CLIP [27] and BERT [4]. These pre-trained text encoders compute text embeddings for natural-language inputs, mapping these inputs to a metric space.

We pre-compute the text embeddings $\theta_d$ of each entry in the dictionary of detected object classes using the pre-trained text encoder, without any fine-tuning. At runtime, given a natural-language prompt specified by a user, FAST-Splat computes the text embedding $\theta_q$ associated with the prompt using the pre-trained text encoder. To identify semantically-relevant objects in the scene, we compute the *relevancy score* [12], which represents the pairwise softmax between the cosine-similarity values computed over the dictionary embeddings $\theta_d$ and a set of text embeddings consisting of the query embedding $\theta_q$ (referred to as the *positive* embedding) and a set of canonical embeddings (referred to as *negative* embeddings).

Given a threshold on the relevancy score, we designate all entries whose relevancy score exceeds this specified threshold as being similar to the query. We map these entries to pixels in a rendered feature map using the probability distribution over all entries computed via the softmax function. In general, we can assign a label to each pixel by taking the maximum probability value; however, alternative techniques can be utilized, including those that provide a multi-hypothesis prediction set, e.g., conformal prediction [34].

**Fine-grained Language Semantics.** In the real-world, identifying only the location of a related object of interest, e.g., in semantic object localization, is generally inadequate. Identifying the precise semantic class of the localized object is of even greater importance, especially for downstream applications where semantic object localization constitutes a single component of a solution pipeline. Prior semantic Gaussian Splatting methods lack this critical capability. For example, existing semantic Gaussian Splatting methods will localize a coffee machine when prompted with the query "tea," without disambiguating the semantic identity of the localized object, which could be misleading to a user. FAST-Splat seeks to address this limitation. The novel distillation technique employed by FAST-Splat endows the algorithm with the capa-

bility to not only identify the precise semantic class of each object, but also to disambiguate the semantic identity of related objects. Given a natural-language prompt, FAST-Splat identifies the related object classes from its dictionary and subsequently, provides the location of these objects in the scene along with their specific semantic label, illustrated in Figure 1. We summarize the architecture in Figure 2.

## 5. Experiments

We benchmark FAST-Splat against the primary, existing semantic Gaussian Splatting methods LangSplat [26] and Feature-3DGS [49]. Unfortunately, FMGS [51], another primary method, does not provide a publicly-available code implementation, preventing us from benchmarking against this method. We utilize the original code implementation provided by the authors of these baseline methods. In addition, we examine the capability of FAST-Splat in disambiguating the semantic identities of objects, e.g., in scenarios with ambiguous user prompts and scenarios with multiple related objects. We examine each method across five scenes: *Kitchen*; *Cleaning*; *Meal*; *Workshop*; and *Plants*. We abbreviate the name of each scene using the initial letter. We utilize YOLO [40] as the closed-set object detector; we provide additional implementation details and further discussion in the Supplementary Material.

**Semantic Gaussian Splatting.** We examine the fidelity of the 3D scene reconstruction achieved by FAST-Splat in comparison to prior work LangSplat [26] and Feature-3DGS [49], beginning with a discussion of the computation time required for data pre-processing as well as the training time for each method, reported in Table 1. FAST-Splat requires at least a 13x faster data pre-processing time compared to the best-competing method Feature-3DGS, and is about 78x faster when compared to LangSplat. Similarly, FAST-Splat achieves about 5x (and up to 8x) faster training times compared to Feature-3DGS, and is about 4x to 6x faster to train compared to LangSplat.

We utilize standard metrics for comparing 3D reconstruction methods, such as the peak-signal-to-noise ratio (PSNR) and Structural Similarity Index Measure (SSIM) for evaluating the reconstruction quality; the training time; the memory

Table 1. The data pre-processing time (prep.) and training time (train) for the semantic Gaussian Splatting methods across five scenes. FAST-Splat outperforms all existing methods.

| Scene | Method | Prep. (min) ↓ | Train (min) ↓ |
|---|---|---|---|
| K. | LangSplat [26] | 128.49 | 58.54 |
| | Feature-3DGS [49] | 15.09 | 73.15 |
| | FAST-Splat [**ours**] | **1.12** | **12.08** |
| C. | LangSplat [26] | 85.53 | 73.32 |
| | Feature-3DGS [49] | 16.49 | 85.10 |
| | FAST-Splat [**ours**] | **1.01** | **12.44** |
| M. | LangSplat [26] | 82.67 | 68.22 |
| | Feature-3DGS [49] | 16.36 | 68.01 |
| | FAST-Splat [**ours**] | **1.06** | **11.71** |
| W. | LangSplat [26] | 103.20 | 66.32 |
| | Feature-3DGS [49] | 15.17 | 65.27 |
| | FAST-Splat [**ours**] | **0.81** | **12.18** |
| P. | LangSplat [26] | 136.99 | 62.61 |
| | Feature-3DGS [49] | 16.05 | 110.86 |
| | FAST-Splat [**ours**] | **1.23** | **13.79** |

Table 2. The rendering frame-rate (Speed), image quality, and memory usage (Memory) achieved by the semantic Gaussian Splatting methods across the five scenes.

| Scene | Method | Speed ↑ (F/s) | SSIM ↑ | PSNR ↑ | Memory ↓ (GB) |
|---|---|---|---|---|---|
| K. | LangSplat [26] | 2.25 | 0.966 | **35.26** | 8.05 |
| | Feature-3DGS [49] | 0.23 | 0.931 | 30.20 | 7.72 |
| | FAST-Splat [**ours**] | **115.23** | **0.971** | 34.24 | **1.02** |
| C. | LangSplat [26] | 2.25 | 0.918 | **32.51** | 9.54 |
| | Feature-3DGS [49] | 0.22 | 0.871 | 28.70 | 7.72 |
| | FAST-Splat [**ours**] | **114.37** | **0.929** | 30.35 | **1.00** |
| M. | LangSplat [26] | 2.24 | 0.933 | **33.37** | 9.54 |
| | Feature-3DGS [49] | 0.22 | 0.845 | 27.83 | 7.86 |
| | FAST-Splat [**ours**] | **127.41** | **0.940** | 31.25 | **0.98** |
| W. | LangSplat [26] | 1.87 | **0.921** | **31.41** | 10.08 |
| | Feature-3DGS [49] | 0.22 | 0.357 | 8.11 | 7.72 |
| | FAST-Splat [**ours**] | **141.48** | 0.917 | 29.22 | **1.04** |
| P. | LangSplat [26] | 6.84 | 0.915 | 28.79 | 7.09 |
| | Feature-3DGS [49] | 0.21 | 0.773 | 25.33 | 7.72 |
| | FAST-Splat [**ours**] | **128.77** | **0.942** | **31.12** | 1.17 |

usage; and the rendering frame-rate. In Table 2, we report the performance of each method on these metrics. FAST-Splat outperforms all other methods on the rendering frame-rate and memory-usage metrics, achieving the highest rendering speeds while requiring the least amount of GPU memory. In four of the five scenes, FAST-Splat achieves the highest SSIM score, outperforming the next best-competing method LangSplat. LangSplat achieves the highest PSNR score in four of the five scenes, followed by FAST-Splat. FAST-Splat achieves the highest PSNR score in the *Plants* scene. Further, FAST-Splat reduces the memory usage in semantic Gaussian Splatting by a factor of more than 3x compared to the best-competing method LangSplat and a factor of more than 6x compared to Feature-3DGS. In essence, FAST-Splat achieves a higher-fidelity scene reconstruction with a lower memory usage.

Lastly, we present the rendering speed of each method. FAST-Splat attains 18x to 75x faster rendering speeds compared to the best-competing method LangSplat. Similarly, FAST-Splat achieves about 550x faster-rendering frame-rates compared to Feature-3DGS, highlighting the superior performance of FAST-Splat. We show rendered images from each of the scenes with the different methods in Figure 3, juxtaposed with the ground-truth RGB image. We render novel views in FAST-Splat.

**Semantic Object Localization.** Now, we evaluate each method on semantic object segmentation from open-vocabulary queries, using the mean intersection-over-union (mIoU) and accuracy metrics. In Table 3, we report the performance of each method on these metrics. FAST-Splat achieves the best performance in two scenes: the *Kitchen* scene and the *Workshop* scene. In the *Plants* scene, Feature-3DGS achieves the highest mIoU score, while LangSplat achieves the highest accuracy. Further, in the *Cleaning* scene, Feature-3DGS achieves the highest mIoU score, while FAST-Splat achieves the highest accuracy. In

all scenes where FAST-Splat does not achieve the best score, FAST-Splat achieves the next-best score, remaining competitive with the best-performing method, which underscores the overall effectiveness of our proposed method. In Figure 4, we show segmentation masks for each of the methods across the five scenes.

Notably, all methods perform poorly in the *Workshop* scene. We believe this observation may be due to the *Workshop* scene being mostly out-of-distribution, considering that the training-data distribution of most vision models and vision-language models do not contain industrial objects, such as those found in workshops. Recall that we utilize the YOLO model for closed-set object detection. The YOLO model is trained using the object categories found in the COCO dataset [19], which does not contain common objects found in a workshop, such as a power drill and screwdriver, which is present in the *Workshop* scene. Moreover, other open-vocabulary object detection methods also struggle to identify industrial objects, e.g., GroundingDINO [22].

**Semantic Disambiguation.** FAST-Splat extends the capabilities of semantic Gaussian Splatting methods beyond semantic object localization to the domain of semantic disambiguation. As noted in the preceding discussion, existing semantic Gaussian Splatting methods cannot resolve semantic ambiguity either arising from the user-provided natural-language query or inherently from the scene composition. In contrast, when provided with an ambiguous natural-language query, FAST-Splat disambiguates the semantic object localization task, by identifying the specific semantic label of related objects in the scene.

We provide some examples in Figure 5. For example, when a user queries for "coffee" in the *Kitchen* scene, FAST-Splat provides the user with the semantic identity of the related object *coffee machine,* along with the semantic segmentation mask and location of the coffee machine,
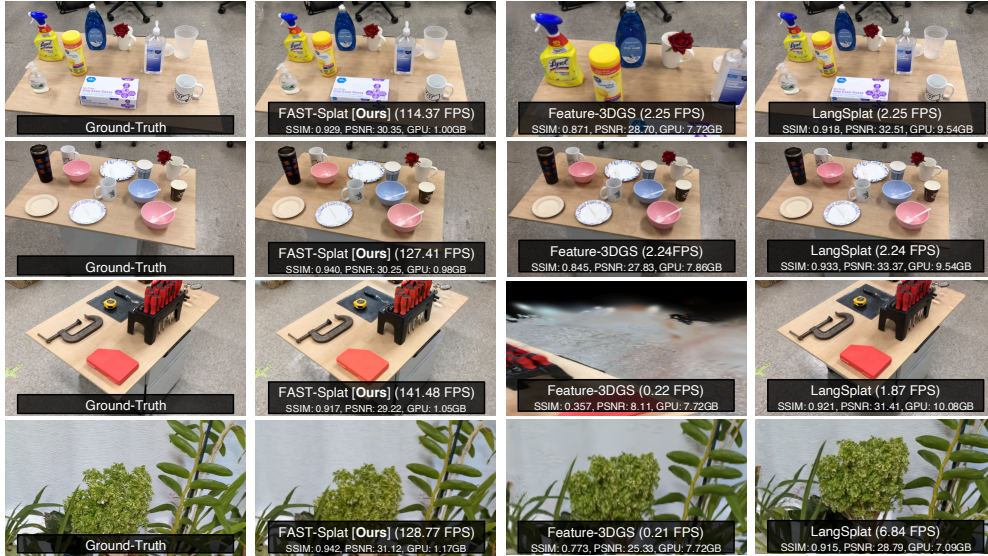
Figure 3. RGB images rendered using each semantic Gaussian Splatting method. FAST-Splat achieves 18x to 75x faster rendering speeds with at least 3x lower memory usage, while achieving competitive (and in some cases, better) reconstruction quality.
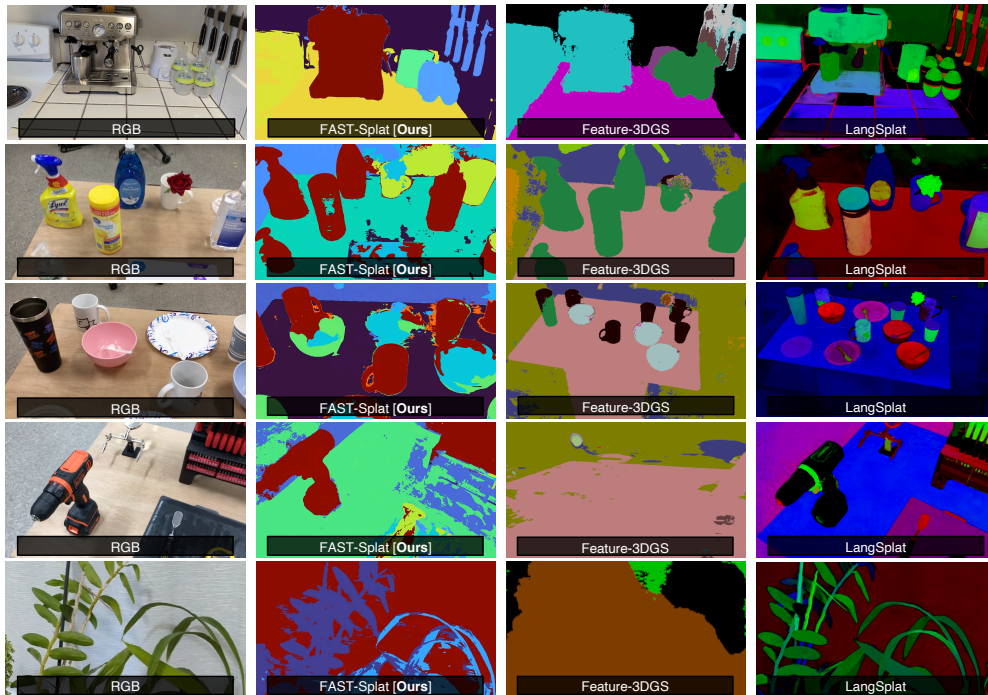


Figure 4. Semantic segmentation using each semantic Gaussian Splatting method. FAST-Splat achieves competitive performance compared to existing semantic Gaussian Splatting method, achieving the highest mIoU and accuracy scores in some scenes while being the next best-performing method in the remaining scenes.

given a threshold for the segmentation task. Likewise, when queried for a "cooking pot," FAST-Splat identifies the semantic class: *kettle*, and the segmentation mask and location for the kettle. Notably, in the *Plants* scene, FAST-Splat resolves the ambiguity between a cup and a vase. When queried with the prompt "cup," FAST-Splat informs the user that, more precisely, a *vase* exists in the scene. Similarly, when asked to localize a "fruit," FAST-Splat clarifies that the scene contains a *pottedplant* and provides the semantic segmentation mask and location.

Figure 5. FAST-Splat resolves language ambiguity in user-provided natural-language queries in semantic object localization. FAST-Splat identifies the specific semantic class of each relevant object, e.g., a *coffee machine* and a *kettle*, when prompted with an ambiguous query, e.g., "coffee" and "cooking pot," respectively. Likewise, FAST-Splat disambiguates between a "cup" and a *vase* and between a "fruit" and a *pottedplant* in the garden-like scene.

Table 3. The mean intersection-over-union (mIoU) and accuracy of the segmentation masks generated by semantic Gaussian Splatting methods across the five scenes.

| Scene | Method | mIoU ↑ | Accuracy ↑ |
|---|---|---|---|
| K. | LangSplat [26] | 0.525 | 0.920 |
| | Feature-3DGS [49] | 0.550 | 0.920 |
| | FAST-Splat [**ours**] | **0.709** | **0.925** |
| C. | LangSplat [26] | 0.552 | 0.850 |
| | Feature-3DGS [49] | **0.570** | 0.751 |
| | FAST-Splat [**ours**] | 0.517 | **0.901** |
| M. | LangSplat [26] | 0.409 | 0.714 |
| | Feature-3DGS [49] | **0.571** | **0.913** |
| | FAST-Splat [**ours**] | 0.525 | 0.904 |
| W. | LangSplat [26] | 0.054 | 0.350 |
| | Feature-3DGS [49] | 0.061 | 0.327 |
| | FAST-Splat [**ours**] | **0.127** | **0.908** |
| P. | LangSplat [26] | 0.562 | **0.900** |
| | Feature-3DGS [49] | **0.641** | 0.872 |
| | FAST-Splat [**ours**] | 0.570 | 0.861 |

**Scene Editing.** Further, like prior semantic Gaussian Splatting methods, FAST-Splat enables scene-editing of Gaussian Splatting scenes. With FAST-Splat, a user can modify the visual properties of objects in the scene, and in addition, can insert, delete, or move objects in the scene. For space considerations, we omit a detailed discussion of this capability, given that it has been discussed in prior work. In Figure 6, we provide an example where the visual properties of the *coffee machine* is modified.
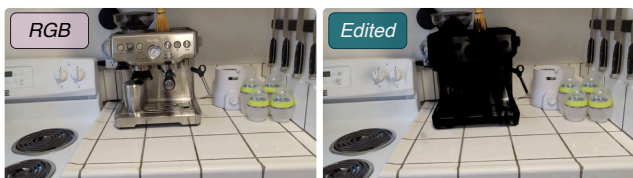


Figure 6. Like existing semantic Gaussian Splatting methods, FAST-Splat enables editing of Gaussian Splatting scenes. Here, we show color-editing of a coffee machine in a kitchen.

# 6. Conclusion

We introduce FAST-Splat, a fast semantic Gaussian Splatting method that resolves ambiguity in semantic object localization. FAST-Splat retains the explicit form of Gaussian Splatting and avoids neural architectures to achieve fast training and rendering speeds compared to existing methods. Further, FAST-Splat builds upon closed-set semantic distillation to enable open-vocabulary semantic Gaussian Splatting. We leverage this formulation to endow FAST-Splat with the capability to resolve semantic ambiguity arising from the natural-language query or the scene composition, in the semantic object localization task, enabling FAST-Splat to address one of the fundamental limitations of existing semantic Gaussian Splatting methods.

# 7. Limitations and Future Work

As with existing semantic Gaussian Splatting methods, the performance of FAST-Splat is highly impacted by the detection performance of the closed-set object detector. Currently, the size of the object categories detected by open-source closed-set object detectors is quite limited. For example, YOLO [40] can only detect 80 objects, which is quite small, considering the diversity of objects in the real-world. This limitation is not imposed by the underlying model architecture. In fact, YOLO can be trained on a much greater diversity of objects. However, the model was trained using the COCO dataset, introducing this limitation. In comparison, existing vision-language models are trained on more diverse datasets, enabling these models to detect a much broader range of objects. Future work will seek to utilize closed-set object detectors trained with a much larger set of object categories.

Our work utilizes an open-vocabulary object detector, e.g., GroundingDINO [22], in addition to a closed-set object detector, which enables FAST-Splat to mitigate the limitations of the closed-set object detector. Essentially, we can always augment the set of detectable objects through natural-

language when querying the open-vocabulary object detector. In fact, we are currently exploring using image-tagging models, e.g., RAM [47], or multi-modal vision-language models, e.g., LLaVa [20] to detect objects for language grounding. Future work will seek to explore this research direction.

## Acknowledgment

## 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. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5855–5864, 2021. 3

[2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020. 2, 4

[3] Boyuan Chen, Fei Xia, Brian Ichter, Kanishka Rao, Keerthana Gopalakrishnan, Michael S Ryoo, Austin Stone, and Daniel Kappler. Open-vocabulary queryable scene representations for real world planning. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11509–11522. IEEE, 2023. 3

[4] Jacob Devlin. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 5, 12

[5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics*, 2019. 2

[6] Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2

[7] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision and language knowledge distillation. *arXiv preprint arXiv:2104.13921*, 2021. 2

[8] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 2

[9] Xu Hu, Yuxi Wang, Lue Fan, Junsong Fan, Junran Peng, Zhen Lei, Qing Li, and Zhaoxiang Zhang. Semantic anything in 3d gaussians. *arXiv preprint arXiv:2401.17857*, 2024. 3, 4

[10] 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. 3

[11] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 3, 4, 12

[12] 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. 2, 3, 4, 5

[13] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014. 4, 12

[14] 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, 3, 4

[15] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and René Ranftl. Language-driven semantic segmentation. *arXiv preprint arXiv:2201.03546*, 2022. 1, 3

[16] Feng Li, Hao Zhang, Huaizhe Xu, Shilong Liu, Lei Zhang, Lionel M Ni, and Heung-Yeung Shum. Mask dino: Towards a unified transformer-based framework for object detection and segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3041–3050, 2023. 3

[17] Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, et al. Grounded language-image pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10965–10975, 2022. 2

[18] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yinan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted clip. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7061–7070, 2023. 3

[19] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 6

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

[21] Quande Liu, Youpeng Wen, Jianhua Han, Chunjing Xu, Hang Xu, and Xiaodan Liang. Open-world semantic segmentation via contrasting and clustering vision-language embedding. In *European Conference on Computer Vision*, pages 275–292. Springer, 2022. 3

[22] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 2, 4, 6, 8

[23] 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. 3

[24] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, et al. Simple open-vocabulary object detection. In *European Conference on Computer Vision*, pages 728–755. Springer, 2022. 1, 2

[25] 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. 3

[26] Minghan Qin, Wanhua Li, Jiawei Zhou, Haoqian Wang, and Hanspeter Pfister. Langsplat: 3d language gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20051–20060, 2024. 2, 3, 4, 5, 6, 8

[27] 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. 1, 2, 3, 5, 12

[28] Adam Rashid, Satvik Sharma, Chung Min Kim, Justin Kerr, Lawrence Yunliang Chen, Angjoo Kanazawa, and Ken Goldberg. Language embedded radiance fields for zero-shot task-oriented grasping. In *7th Annual Conference on Robot Learning*, 2023. 2, 3

[29] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 1, 3, 4

[30] N Reimers. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*, 2019. 3

[31] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *IEEE transactions on pattern analysis and machine intelligence*, 39(6):1137–1149, 2016. 2

[32] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, et al. Grounded sam: Assembling open-world models for diverse visual tasks. *arXiv preprint arXiv:2401.14159*, 2024. 3

[33] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 12

[34] Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. *Journal of Machine Learning Research*, 9(3), 2008. 5

[35] Nur Muhammad Mahi Shafiullah, Chris Paxton, Lerrel Pinto, Soumith Chintala, and Arthur Szlam. Clip-fields: Weakly supervised semantic fields for robotic memory. *arXiv preprint arXiv:2210.05663*, 2022. 2, 3

[36] 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. 2, 3

[37] Ola Shorinwa, Johnathan Tucker, Aliyah Smith, Aiden Swann, Timothy Chen, Roya Firoozi, Monroe Kennedy III, and Mac Schwager. Splat-MOVER: Multi-stage, open-vocabulary robotic manipulation via editable gaussian splatting. *arXiv preprint arXiv:2405.04378*, 2024. 3

[38] Matthew Tancik, Ethan Weber, Evonne Ng, Ruilong Li, Brent Yi, Terrance Wang, Alexander Kristoffersen, Jake Austin, Kamyar Salahi, Abhik Ahuja, et al. Nerfstudio: A modular framework for neural radiance field development. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–12, 2023. 12

[39] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Clip-nerf: Text-and-image driven manipulation of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3835–3844, 2022. 2, 3, 4

[40] Chien-Yao Wang, I-Hau Yeh, and Hong-Yuan Mark Liao. Yolov9: Learning what you want to learn using programmable gradient information. In *European Conference on Computer Vision*, pages 1–21. Springer, 2025. 4, 5, 8, 12

[41] Weiyao Wang, Matt Feiszli, Heng Wang, Jitendra Malik, and Du Tran. Open-world instance segmentation: Exploiting pseudo ground truth from learned pairwise affinity. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4422–4432, 2022. 3

[42] Zhuyu Yao, Jiangbo Ai, Boxun Li, and Chi Zhang. Efficient detr: improving end-to-end object detector with dense prior. *arXiv preprint arXiv:2104.01318*, 2021. 4

[43] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917*, 2022. 1

[44] Alireza Zareian, Kevin Dela Rosa, Derek Hao Hu, and Shih-Fu Chang. Open-vocabulary object detection using captions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14393–14402, 2021. 2

[45] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung Shum. Dino: Detr with improved denoising anchor boxes for end-to-end object detection. *arXiv preprint arXiv:2203.03605*, 2022. 2, 3

[46] Hao Zhang, Feng Li, Xueyan Zou, Shilong Liu, Chunyuan Li, Jianwei Yang, and Lei Zhang. A simple framework for open-vocabulary segmentation and detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1020–1031, 2023. 3

[47] Youcai Zhang, Xinyu Huang, Jinyu Ma, Zhaoyang Li, Zhaochuan Luo, Yanchun Xie, Yuzhuo Qin, Tong Luo, Yaqian Li, Shilong Liu, et al. Recognize anything: A strong image tagging model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1724–1732, 2024. 9

[48] Yuhang Zheng, Xiangyu Chen, Yupeng Zheng, Songen Gu, Runyi Yang, Bu Jin, Pengfei Li, Chengliang Zhong, Zengmao Wang, Lina Liu, et al. Gaussiangrasper: 3d language gaussian splatting for open-vocabulary robotic grasping. *arXiv preprint arXiv:2403.09637*, 2024. 3

[49] 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, 3, 4, 5, 6, 8

[50] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In *European Conference on Computer Vision*, pages 350–368. Springer, 2022. 1

[51] Xingxing Zuo, Pouya Samangouei, Yunwen Zhou, Yan Di, and Mingyang Li. Fmgs: Foundation model embedded 3d gaussian splatting for holistic 3d scene understanding. *International Journal of Computer Vision*, pages 1–17, 2024. 3, 4, 5

# Appendix A. Implementation Details

To implement FAST-Splat, we build off Splatfacto in Nerfstudio [38], and utilize the text encoder of the OpenAI's CLIP ResNet model *RN50x64* [27] to compute the text embeddings. However, other text encoders can also be used, e.g., other CLIP models or BERT [4]. We train FAST-Splat using Stochastic Gradient Descent with momentum, via the Adam optimizer [13].

We collected the training data by taking videos of the scene using a smartphone camera. We utilized structure-from-motion, e.g., COLMAP [33], for computing the pose of the camera for each frame used in training the Gaussian Splatting model. FAST-Splat is trained on an NVIDIA GeForce RTX 3090 GPU with 24GB VRAM for 30000 iterations, using the hyperparameters specified in [38] in its Splafacto model. We train the baseline methods LangSplat and Feature-3DGS using the original code implementation provided by their respective authors on an NVIDIA GeForce RTX 3090 Ti GPU with 24GB VRAM, giving LangSplat and Feature-3DGS an advantage on hardware, since the 3090 Ti has 10752 CUDA cores compared to the 10496 CUDA cores on the 3090, among other improved hardware specifications. To train LangSplat, we follow the procedure provided by its authors, which involves pre-training the Gaussian Splatting Scene using the code implementation provided in [11], prior to training the semantic components. We do not include the structure-from-motion processing times in any of the results reported in the experiments, given that this procedure is required universally by all methods.

# Appendix B. Semantic Object Localization

Here, we provide additional results, demonstrating the photorealistic novel-view synthesis capability of FAST-Splat, including its semantic segmentation masks. In Figure 7, we show the semantic localization of a "wall," "countertop," "kettle," and "cooking spoon." In each case, FAST-Splat successfully localizes the pertinent object. The colors in each figure vary with the similarity of the object to the natural-language query. We utilize the "Turbo" colormap, with an increase in the intensity of the red color channel indicating an *increase* in the similarity of an object to the query and an increase in the intensity of the blue color channel indicating a *decrease* in the relevance of an object to the query, in general. Moreover, the RGB images rendered by FAST-Splat appear photorealistic, displayed in the top row in Figure 7, all rendered at novel views.

# Appendix C. Semantic Disambiguation

As discussed earlier in this work, FAST-Splat enables semantic disambiguation, e.g., in semantic object localization. Not only does FAST-Splat provide the semantic identity of the most relevant object, FAST-Splat also identifies other pertinent objects, providing the semantic label of these objects along with their similarity to the natural-language prompt. In Figure 8, we demonstrate the semantic disambiguation capabilities of FAST-Splat. In the top row, when queried for a "fork," FAST-Splat identifies the *spoon* as the object that is most similar to the query. In addition, FAST-Splat localizes the *knife*, identifying it as being relevant, although less similar, compared to the spoon. Likewise, in the bottom row, when prompted with the ambiguous query "water," FAST-Splat localizes the *kettle*, noting its lower similarity compared to a *coffee machine*, which FAST-Splat also localizes in the scene. Lastly, FAST-Splat identifies the *bottle* as the most similar object to the query, resolving the ambiguous prompt.

# Appendix D. Limitations

We discuss the limitations of FAST-Splat further. The performance of FAST-Splat may degrade in scenes that are out of the distribution of the training data for the closed-set or open-set object detector used in extracting the object classes. For example, the workshop scene lies primarily outside the distribution of the training data of YOLO [40]. Consequently, FAST-Splat performs relatively poorly in this scene. Figure 9 shows the semantic localization results for a "spatula" and "tape." FAST-Splat fails to localize these objects accurately. However, this limitation can be mitigated by training the closed-set detector on a broader set of object classes; augmenting the query passed into the open-vocabulary object detector; or training the open-vocabulary object data on more data, depending on the primary cause.
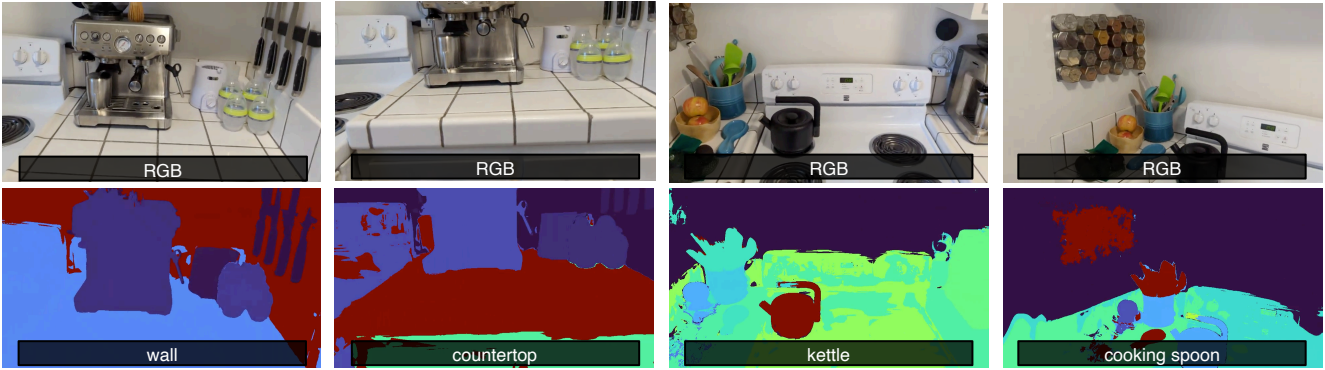
Figure 7. FAST-Splat enables semantic object localization wth open-vocabulary user prompts. Here, we provide segmentation results for a "wall," "countertop," "kettle," and "cooking spoon."



Figure 8. FAST-Splat resolves semantic ambiguity in object localization from natural-language queries. In the top row, although the queried object "fork" does not exist in the scene, FAST-Splat localizes relevant objects in the scene, and more importantly, provides the semantic label of each of these objects. Here, FAST-Splat identifies a *spoon* and a *knife* and notes that the spoon (not the knife) is more similar to the fork. Likewise, in the bottom row, given the ambiguous prompt "water," FAST-Splat localizes a *kettle*, *coffee machine*, and *bottle*, providing their semantic object classes along with their relative similarity to the prompt.



Figure 9. The performance of FAST-Splat degrades in scenes that are out-of-distribution for the object detector. Here, FAST-Splat fails to generate accurate segmentation masks in the *Workshop* scene, which contains objects that are not found in the training dataset of many object detectors.