

3D Feature Distillation with Object-Centric Priors

Georgios Tzifas

Department of Artificial Intelligence
University of Groningen, the Netherlands
g.t.tzifas@rug.nl

Yucheng Xu

School of Informatics
University of Edinburgh, United Kingdom
Yucheng.Xu@ed.ac.uk

Zhibin Li

Department of Computer Science
University College London, United Kingdom
alex.li@ucl.ac.uk

Hamidreza Kasaei

Department of Artificial Intelligence
University of Groningen, the Netherlands
hamidreza.kasaei@rug.nl

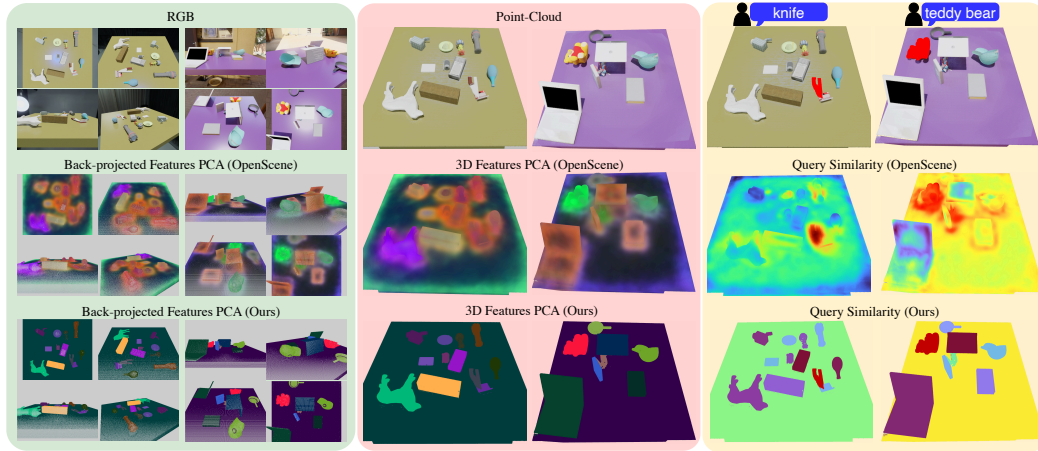


Figure 1: Visualization of 3D features (*middle*), back-projected 2D features (*left*) and query similarity heatmaps (*right*), for OpenScene and our DRO-CLIP. OpenScene fuses pixel-wise 2D features with average pooling, leading to grounding failures and fuzzy object boundaries. Our method tackles such issues using object-centric priors to fuse object-level 2D features in 3D instance masks with semantics-informed view selection.

Abstract: Grounding natural language to the physical world is a ubiquitous topic with a wide range of applications in computer vision and robotics. Recently, 2D vision-language models such as CLIP have been widely popularized, due to their impressive capabilities for open-vocabulary grounding in 2D images. Subsequent works aim to elevate 2D CLIP features to 3D via feature distillation, but either learn neural fields that are scene-specific and hence lack generalization, or focus on indoor room scan data that require access to multiple camera views, which is not practical in robot manipulation scenarios. Additionally, related methods typically fuse features at pixel-level and assume that all camera views are equally informative. In this work, we show that this approach leads to sub-optimal 3D features, both in terms of grounding accuracy, as well as segmentation crispness. To alleviate this, we propose a multi-view feature fusion strategy that employs object-centric priors to eliminate uninformative views based on semantic information, and fuse features at object-level via instance segmentation masks. To distill our object-centric 3D features, we generate a large-scale synthetic multi-view dataset of cluttered tabletop scenes, spawning 15k scenes from over 3300 unique object instances, which we make publicly available. We show that our method reconstructs 3D CLIP features with improved grounding capacity and spatial consistency, while doing so from single-view RGB-D, thus departing from the assumption of multiple camera views at test time. Finally, we show that our approach can generalize to novel tabletop domains and be re-purposed for 3D instance segmentation without fine-tuning, and demonstrate its utility for language-guided robotic grasping in clutter.

1 Introduction

Language grounding in 3D environments plays a crucial role in realizing intelligent systems that can interact naturally with the physical world. In the robotics field, being able to precisely segment desired objects in 3D based on open language queries (object semantics, visual attributes, affordances, etc.) can serve as a powerful proxy for enabling open-ended robot manipulation. As a result, research focus on 3D segmentation methods has seen growth in recent years [1, 2, 3, 4, 5, 6]. However, related methods fall in the closed-vocabulary regime, where only a fixed list of classes can be used as queries. Inspired by the success of open-vocabulary 2D methods [7, 8, 9, 10], recent efforts elevate 2D representations from pretrained image models [7, 11] to 3D via distillation pipelines [12, 13, 14, 15, 16, 17, 18, 19]. In this work, we identify several limitations of existing distillation approaches. On the one hand, field-based methods [13, 20, 16, 17, 18] offer continuous 3D feature fields, but require to be trained online in specific scenes and hence cannot generalize to novel object instances and compositions, they require a few minutes to train, and need to collect multiple camera views before training, all of which hinder their real-time applicability. On the other hand, original 3D feature distillation methods and follow up work [12, 14, 21] use room scan datasets [22, 23] to distill 2D features fused from multiple views with point-cloud encoders. The distilled features maintain the open-set generalizability of the pretrained model, therefore granting such methods applicable in novel scenes with open vocabularies. However, such approaches assume that 2D features from all views are equally informative, which is not the case in natural indoors scenes, where due to partial visibility and clutter, certain views will lead to noisy representations. 2D features are also typically fused point-wise from ViT patches [9, 10, 8] or multi-scale crops [13, 6], therefore leading to the so called “patchyness” issue [24] (see Fig. 1), where features computed in patches / crops that involve multiple objects lead to fuzzy segmentation boundaries. The latter issue is especially impactful in robot manipulation, where precise 3D segmentation is vital for specifying robust actuation goals.

To address such limitations, in this work, we revisit $2D \rightarrow 3D$ feature distillation with point-cloud encoders, but revise the multi-view feature fusion strategy to enhance the quality of the target 3D features. In particular, we inject both *semantic* and *spatial* object-centric priors into the fusion strategy, in three ways: (i) We obtain object-level 2D features by isolating object instances in each camera view from their 2D segmentation masks, (ii) we fuse features only at corresponding 3D object regions using their 3D segmentation masks, (iii) we leverage dense object-level semantic information to devise an informativeness metric, which is used to weight the contribution of views and eliminate uninformative ones. Extensive ablation studies demonstrate the advantages of our proposed object-centric fusion strategy compared to vanilla approaches. To train our method, we require a large-scale cluttered indoors dataset with dense number of views per scene, which is currently not existent. To that end, we build MV-TOD (**M**ulti-**V**iew **T**abletop **O**bjects **D**ataset), consisting of $\sim 15k$ Blender scenes from more than $3.3k$ unique 3D object models, for which we provide 73 views per scene with 360° coverage, further equipped with 2D/3D segmentations, 6-DoF grasps and semantic object-level annotations. We use MV-TOD to distill the object-centric 3D CLIP [7] features acquired via our fusion strategy into a 3D representation, which we call DROP-CLIP (**D**istilled **R**epresentations with **O**bject-centric **P**riors from CLIP). Our 3D encoder operates in partial point-clouds from a single RGB-D view, thus departing from the requirement of multiple camera images at test time, while offering real-time inference capabilities. By imposing the same 3D features as distillation targets for a large number of diverse views, we encourage DROP-CLIP to learn a view-invariant 3D representation. We demonstrate that our learned 3D features surpass previous 3D open-vocabulary approaches in semantic and referring segmentation tasks in MV-TOD, both in terms of grounding accuracy and segmentation crispness, while significantly outperforming previous 2D approaches in the single-view setting. Further, we show that they can be leveraged zero-shot in novel tabletop datasets that contain real-world scenes with unseen objects and new vocabulary, as well as be used out-of-the-box for 3D instance segmentation tasks, performing competitively with established segmentation approaches without fine-tuning.

In summary, our contributions are fourfold: (i) we release MV-TOD, a large-scale synthetic dataset of household objects in cluttered tabletop scenarios, featuring dense multi-view coverage and semantic/mask/grasp annotations, (ii) we identify limitations of current multi-view feature fusion

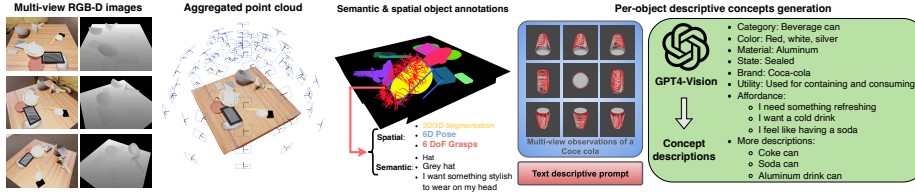


Figure 2: **MV-TOD Overview:** Example generated scene, source multi-view RGB=D images and scene annotations (*left*). Automatic semantic annotation generation with VLMS (*right*).

approaches and illustrate how to overcome them by leveraging object-centric priors, (iii) we release DROP-CLIP, a 3D model that reconstructs view-independent 3D CLIP features from single-view, and (iv) we conduct extensive ablation studies, comparative experiments and robot demonstrations to showcase the effectiveness of the proposed method in terms of 3D segmentation performance, generalization to novel domains and tasks, and applicability in robot manipulation scenarios.

2 Multi-View Tabletop Objects Dataset

Existing 3D datasets mainly focus on indoor scenes in room layouts [33, 22, 26] and related annotations typically cover closed-set object categories (e.g. furniture) [1, 2, 27, 34, 28], which are not practical for robot manipulation tasks, where cluttered tabletop scenarios and open-vocabulary language are of key importance. Alternatively, recent

Dataset	Layout	Multi View	Clutter	Vision Data	Ref.Expr. Annot.	Grasp Annot.	Num.Obj. Categories	Num. Scenes	Num. Expr.	Obj-lvl Semantics
ScanNet [22]	indoor	✓	-	RGB-D,3D	✗	✗	17	800	-	✗
S3DIS [25]	indoor	✓	-	RGB-D,3D	✗	✗	13	6	-	✗
Replica [26]	indoor	✓	-	RGB-D,3D	✗	✗	88	-	-	✗
STPLS3D [25]	outdoor	✓	-	3D	✗	✗	12	18	-	✓
ScanRefer [1]	indoor	✓	✗	RGB-D,3D	2D/3D mask	✗	18	800	51.5k	✗
ReferIt-3D [2]	indoor	✓	✗	RGB-D,3D	2D/3D mask	✗	18	707	125.5k	✗
ReferIt-RGBD [27]	indoor	✓	✗	RGB-D	2D box	✗	-	7.6k	38.4k	✗
SunSpot [28]	indoor	✓	✓	RGB-D	2D box	✗	38	1.9k	7.0k	✗
GraspNet [29]	tabletop	✗	✓	3D	✗	6-DoF	88	190	-	✗
REGRAD [30]	tabletop	✗	✓	RGB-D,3D	✗	6-DoF	55	47k	-	✗
OCID-VLG [31]	tabletop	✗	✓	RGB-D,3D	2D mask	4-DoF	31	1.7k	89.6k	template
Grasp-Anything [32]	tabletop	✗	✗	RGB	2D mask	4-DoF	236	1M	-	open
MV-TOD (ours)	tabletop	✓	✓	RGB-D,3D	3D mask	6-DoF	149	15k	671.2k	open

Table 1: Comparisons between MV-TOD and existing datasets. Table 1 summarizes key differences between MV-TOD and existing grounding / grasping datasets.

MV-TOD contains approximately 15k scenes generated in Blender [36], comprising of 3379 unique object models, 99 collected by us and the rest filtered from ShapeNet-Sem model set [37]. The dataset enumerates 149 object categories featuring typical household objects (kitchenware, food, electronics etc.), each of which includes multiple instances that vary in fine-grained details such as color, texture, shape etc. For each object instance, we leverage modern vision-language models such as GPT-4-Vision [38] to generate textual annotations referring to various object attributes, including category, color, material, state, utility, brand, etc., spawning over 670k unique referring instance queries. We refer the reader to Appendix A.1 for details on object statistics and scene generation implementation. For each scene, we provide 73 uniformly distributed views, 2D / 3D instance segmentation masks, 6D object poses, as well as a set of referring expressions sampled from the object-level semantic annotations. Additionally, we provide collision-free 6-DoF grasp poses for each scene object, originating from the ACRONYM dataset [35]. In this paper, we leverage the dense multi-view coverage of MV-TOD for 2D → 3D feature distillation. However, given the breadth of labels in MV-TOD, we believe it can serve as a resource for several 3D vision and robotics

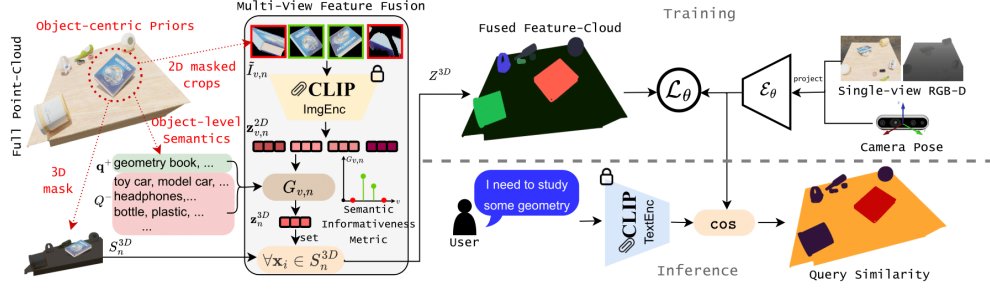


Figure 3: **Method Overview:** Given a 3D scene and multiple camera views, we employ three object-centric priors (*in red*) for multi-view feature fusion: (i) extract CLIP features from 2D masked object crops, (ii) use semantic annotations to fuse 2D features across views, (iii) apply the fused feature on all points in the object’s 3D mask. The fused feature-cloud is distilled with a single-view posed RGB-D encoder and cosine distance loss. During inference, we compute point-wise cosine similarity scores in CLIP space (higher similarity towards red).

downstream tasks, including instance segmentation, 6D pose estimation and 6-DoF grasp synthesis. To the best of our knowledge, MV-TOD is the first dataset to combine 3D cluttered scenes with multi-view images, open-vocabulary language and 6-DoF grasp annotations.

3 Distilled Representations with Object-Centric Priors

Our goal is to distill multi-view 2D CLIP features into a 3D representation, while employing an object-centric feature fusion strategy to ensure high quality 3D features. Our overall pipeline is illustrated in Fig 3. We first introduce traditional multi-view feature fusion (Sec. 3.1), present our variant with object-centric priors (Sec. 3.2), discuss feature distillation training (Sec. 3.3) and describe how to perform inference for downstream open-vocabulary 3D grounding tasks (Sec. 3.4).

3.1 Multi-view Feature Fusion

We assume access to a dataset of 3D scenes, where each scene is represented through a set of \mathcal{V} posed RGB-D views $\{I_v \in \mathbb{R}^{H \times W \times 3}, D_v \in \mathbb{R}^{H \times W}, T_v \in \mathbb{R}^{4 \times 4}\}_{v=1}^{\mathcal{V}}$, with $H \times W$ denoting the image resolution, \mathcal{V} the total number of views, and T_v the transformation matrix from each camera’s viewpoint v with respect to a global reference frame, such as the center of the tabletop. A projection matrix K_v representing each camera’s intrinsic parameters is also given. For each scene we reconstruct the full point-cloud $P \in \mathbb{R}^{M \times 3}$ by aggregating all depth images D_v , after projecting them to 3D with the camera intrinsics K_v and transforming to world frame with T_v^{-1} . To remove redundant points, we voxelize the aggregated point-cloud with a fixed voxel size resolution d^3 , resulting in M total points. Our goal is to obtain a feature-cloud $Z^{3D} \in \mathbb{R}^{M \times C}$, where C is the dimension of the representations provided by the pretrained image model, fused across all views.

2D feature extraction We pass each RGB view to a pretrained image model $f^{2D} : \mathbb{R}^{H \times W \times 3} \rightarrow \mathbb{R}^{H \times W \times C}$ to obtain pixel-level features $Z_v^{2D} = f^{2D}(I_v)$. Any ViT-based vision foundation model (e.g. DINO-v2 [11]) can be chosen, but we focus on CLIP [7], since we want our 3D representation to be co-embedded with language, as to enable open-vocabulary grounding. However, vanilla CLIP features are restrained to image-level, whereas we require dense pixel-level features to perform multi-view fusion. To obtain pixel-wise features, previous works explore fine-tuned CLIP models [12, 15] such as OpenSeg [9] or LSeg [10], multi-scale crops from anchored points in the image frame [13, 6, 21] or MaskCLIP [8, 16], which provides patch-level text-aligned features from CLIP’s ViT encoder without additional training. All approaches are compatible with our framework (ablations in Sec. 4.1).

2D-3D correspondence Given the i -th point in P , $\mathbf{x}_i = (x, y, z)$, $i = 1, \dots, M$, we first back-project to each camera view v using: $\tilde{\mathbf{u}}_{v,i} \doteq \mathcal{M}_v(\mathbf{x}_i) = K_v \cdot T_v \cdot \tilde{\mathbf{x}}_i$, where $\tilde{\mathbf{u}} = (u_x, u_y, u_z)^T$ and $\tilde{\mathbf{x}} = (x, y, z, 1)^T$ homogeneous coordinates in 2D camera frame and 3D world frame respectively,

and $\mathbf{u} = (u_x, u_y)^T$. The 2D feature for each back-projected point $\mathbf{z}_{v,i}^{2D} \in \mathbb{R}^C$ is then given by:

$$\mathbf{z}_{v,i}^{2D} = f^{2D}(I_v(\mathbf{u}_{v,i})) = f^{2D}(I_v(\mathcal{M}_v(\mathbf{x}_i))) \quad (1)$$

For each view, we eliminate points that fall outside of a camera view’s FOV by considering only the pixels: $\{\tilde{\mathbf{u}}_v = (u_x, u_y, u_z)^T \in \mathcal{M}_v(P) \mid u_z \neq 0, u_x/u_z \in [0, W), u_y/u_z \in [0, H)\}$. It is further important to maintain only points that are visible from each camera view, as a point might lie within the camera’s FOV but in practise be occluded by a foreground object. To eliminate such points, we follow [12, 6] and compare the back-projected z coordinate u_z with the sensor depth reading $D_v(u_x, u_y)$. We maintain only points that satisfy: $|u_z - D_v(u_x, u_y)| \leq c_{thr}$, where c_{thr} a fixed hyper-parameter. We compose the FOV and occlusion filtering to obtain a *visibility map* $\Lambda_{v,i} \in \{0, 1\}^{V \times M}$, which determines whether point i is visible from view v .

Fusing point-wise features Obtaining a 3D feature for each point $i = 1, \dots, M$ is achieved by fusing back-projected 2D features $\mathbf{z}_{v,i}^{2D}$ with weighted-average pooling:

$$\mathbf{z}_i^{3D} = \frac{\sum_{v=1}^V \mathbf{z}_{v,i}^{2D} \cdot \omega_{v,i}}{\sum_{v=1}^V \omega_{v,i}} \quad (2)$$

where $\omega_{v,i} \in \mathbb{R}$ a scalar weight that represents the importance of view v for point i . In practise, previous works consider $\omega_{v,i} = \Lambda_{v,i}$ [12], a binary weight for the visibility of each point. In essence, this method assumes that all views are equally informative for each point, as long as the point is visible from that view.

We suggest that naively average pooling 2D features for each point leads to sub-optimal 3D features, as noisy, uninformative views contribute equally, therefore “polluting” the overall representation. In our work we propose to decompose $\omega_{v,i} = \Lambda_{v,i} \cdot G_{v,i}$, where $G_{v,i} \in \mathbb{R}^{V \times M}$ an informativeness weight that measures the importance of each view for each point. In the next subsection, we describe how to use text data to dynamically compute an informativeness weight for each view based on *semantic* object-level information, as well as how to perform object-wise instead of point-wise fusion.

3.2 Employing Object-Centric Priors

Let $\{S_v^{2D} \in \{0, 1\}^{N \times H \times W}\}_{v=1}^V$ be view-aligned 2D instance-wise segmentation masks for each scene, where N the total number of scene objects, provided from the training dataset. We aggregate the 2D masks to obtain $S^{3D} \in \{1, \dots, N\}^M$, such that for each point i we can retrieve the corresponding object instance $n_i = S_i^{3D}$.

Semantic informativeness metric Let $\mathcal{Q} = \{Q_k\}_{k=1}^K$, $Q_k \in \mathbb{R}^{N_k \times C}$ be a set of object-specific textual prompts, where K the number of dataset object instances and N_k the number of prompts for object k . We use CLIP’s text encoder to embed the textual prompts in \mathbb{R}^C and average them to obtain an object-specific prompt $\mathbf{q}_k = 1/N_k \cdot \sum_{j=1}^{N_k} Q_{k,j}$. For each scene, we map each object instance $n \in [1, N]$ to its positive prompt \mathbf{q}_n^+ , as well as a set $Q_n^- \doteq \mathcal{Q} - \{\mathbf{q}_n^+\}$ of negative prompts corresponding to all other instances. We define our *semantic informativeness metric* as:

$$G_{v,i} = \cos(\mathbf{z}_{v,i}^{2D}, \mathbf{q}_{n_i}^+) - \max_{\mathbf{q} \sim Q_n^-} \cos(\mathbf{z}_{v,i}^{2D}, \mathbf{q}) \quad (3)$$

Intuitively, we want a 2D feature from view v to contribute to the overall 3D feature of point i according to how much its similarity with the correct object instance is higher than the maximum similarity to any of the negative object instances, hence offering a proxy for semantic informativeness. We clip this weight to 0 to eliminate views that don’t satisfy the condition $G_{v,i} \geq 0$. Plugging in our metric in equation (2) already provides improvements over vanilla average pooling (see Sec. 4.1), however, does not deal with 3D spatial consistency, for which we employ our spatial priors below.

Object-level 2D CLIP features For obtaining object-level 2D CLIP features, we isolate the pixels for each object n from each view v from $S_{v,n}^{2D}$ and crop a bounding box around the mask from I_v : $\mathbf{z}_{v,n}^{2D} = f_{cls}^{2D}(\text{cropmask}(I_v, S_{v,n}^{2D}))$ (see Appendix A.3 for ablations in CLIP visual prompts). Here

we use $f_{cls}^{2D} : \mathbb{R}^{h_n \times w_n \times 3} \rightarrow \mathbb{R}^C$, i.e., only the [CLS] feature of CLIP’s ViT encoder, to represent an object crop of size $h_n \times w_n$. We can now define our metric from equation (3) also at object-level:

$$G_{v,n} = \cos(\mathbf{z}_{v,n}^{2D}, \mathbf{q}_n^+) - \max_{\mathbf{q} \sim Q_n^-} \cos(\mathbf{z}_{v,n}^{2D}, \mathbf{q}) \quad (4)$$

where $G_{v,n} \in \mathbb{R}^{\mathcal{V} \times \mathcal{N}}$ now represents the semantic informativeness of view v for object instance n .

Fusing object-wise features A 3D object-level feature can be obtained by fusing 2D object-level features across views similar to equation (2):

$$\mathbf{z}_n^{3D} = \frac{\sum_{v=1}^{\mathcal{V}} \mathbf{z}_{v,n}^{2D} \cdot \omega_{v,n}}{\sum_{v=1}^{\mathcal{V}} \omega_{v,n}} = \frac{\sum_{v=1}^{\mathcal{V}} \mathbf{z}_{v,n}^{2D} \cdot \Lambda_{v,n} \cdot G_{v,n}}{\sum_{v=1}^{\mathcal{V}} \Lambda_{v,n} \cdot G_{v,n}} \quad (5)$$

where each view is weighted by its semantic informativeness metric $G_{v,n}$, as well as optionally a visibility metric $\Lambda_{v,n} = \sum S_{v,n}^{2D}$ that measures the number of pixels from n -th object’s mask that are visible from view v [6]. We finally reconstruct the full feature-cloud $Z^{3D} \in \mathbb{R}^{M \times C}$ by equating each point’s feature to its corresponding 3D object-level one via: $\mathbf{z}_i^{3D} = \mathbf{z}_{n_i}^{3D}$, $n_i = S_i^{3D}$.

3.3 View-Independent Feature Distillation

Even though the above feature-cloud Z^{3D} could be directly used for open-vocabulary grounding in 3D, its construction is computationally intensive and requires a lot of expensive resources, such as access to multiple camera views, view-aligned 2D instance segmentation masks, as well as textual prompts to compute informativeness metrics. Such utilities are rarely available in open-ended scenarios, especially in robotic applications, where usually only single-view RGB-D images from sensors mounted on the robot are provided. To tackle this, we wish to distill all the above knowledge from the feature-cloud Z^{3D} with an encoder network that receives only a partial point-cloud from single-view posed RGB-D. Hence, the only assumption that we make during inference is access to camera intrinsic and extrinsic parameters, which is a mild requirement in most robotic pipelines.

In particular, given a partial colored point-cloud from view v : $P_v \in \mathbb{R}^{M_v \times 6}$, we train an encoder $\mathcal{E}_\theta : \mathbb{R}^{M_v \times 6} \rightarrow \mathbb{R}^{M_v \times C}$ such that $\mathcal{E}_\theta(P_v) = Z^{3D}$. Notice that the distillation target Z^{3D} is independent of view v . Following [12, 15] we use cosine distance loss:

$$\mathcal{L}(\theta) = 1 - \cos(\mathcal{E}_\theta(P_v), Z^{3D}) \quad (6)$$

See Appendix A.2 for training implementation details. With such a setup, we can obtain 3D features that: (i) are co-embedded in CLIP text space, so they can be leveraged for 3D segmentation tasks from open-vocabulary queries, (ii) are ensured to be optimally informative per object, due to the usage of the semantic informativeness metric to compute Z^{3D} , (iii) maintain 3D spatial consistency in object boundaries, due to performing object-wise instead of point-wise fusion when computing Z^{3D} , and (iv) are encouraged to be view-independent, as the same features Z^{3D} are utilized as distillation targets regardless of the input view v . Importantly, no labels, prompts, or segmentation masks are needed at test-time to reproduce the fused feature-cloud, while obtaining it amounts to a single forward pass of our 3D encoder, hence offering real-time performance.

3.4 Open-Vocabulary 3D Segmentation

Given a predicted feature-cloud $\hat{Z}^{3D} = \mathcal{E}_\theta(P_v)$, we can perform 3D grounding tasks from open-vocabularies by computing cosine similarities between CLIP text embeddings and \hat{Z}^{3D} .

Semantic segmentation In this task, the queries correspond to an open-set of textual prompts $Q = \{\mathbf{q}_k\}_{k=1}^{\mathcal{K}}$ describing \mathcal{K} semantic classes. A class for each point $\hat{Y} \in \{1, \dots, \mathcal{K}\}^M$ is given by: $\hat{Y} = \operatorname{argmax}_k \cos(\hat{Z}^{3D}, \mathbf{q}_k)$.

Referring segmentation Here the user provides an open-vocabulary query \mathbf{q}^+ referring to a particular object instance, and optionally a set of negative prompts $Q^- \in \mathbb{R}^{N^- \times C}$, which in practise can be initialized from an open-set as above or with canonical phrases (e.g. ‘object’, ‘thing’ etc.) [13].

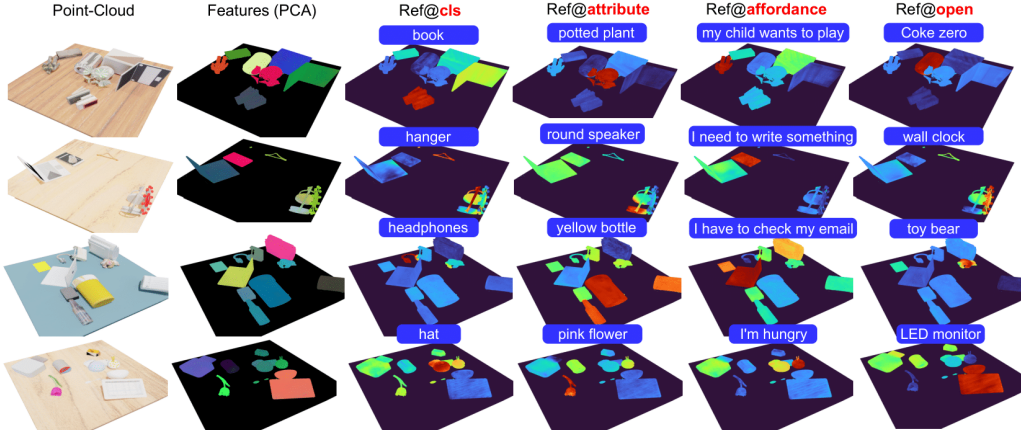


Figure 4: **Open-Vocabulary 3D Referring Segmentation in MV-TOD.** Examples of learned 3D features and grounding heatmaps from open-ended language queries (class names, attributes, user affordances, and open instance-specific concepts) in scenes from MV-TOD dataset. Points are colored based on their query similarity (higher towards red). We note that table points are excluded from similarity computation in our visualizations.

Similarity scores are converted to probabilities: $\mathcal{P} = \text{softmax}\left(\frac{1}{\gamma} \cdot \cos(\hat{Z}^{3D}, [\mathbf{q}^+, \mathbf{Q}^-]^T)\right)$, where γ a temperature hyper-parameter and $\mathcal{P} = [\rho^+, \mathcal{P}^-]$ probabilities of positive matching $\rho^+ \in \mathbb{R}^M$ and negative matching $\mathcal{P}^- = [\rho_1^-, \dots, \rho_{N^-}^-] \in \mathbb{R}^{M \times N^-}$ respectively. The final 3D segmentation is given by $\hat{S}_i = (\rho_i^+ > \max_j \mathcal{P}_{i,j}^-)$, or by thresholding ρ^+ with a fixed threshold s_{thr} (see ablations in Appendix A.3)

Instance segmentation Since our encoder has been distilled with the aid of instance-wise segmentation masks, the obtained features can be utilized out-of-the-box for 3D instance segmentation tasks. We demonstrate that with a simple clustering algorithm over \hat{Z}^{3D} we can obtain 3D instance segmentation masks for cluttered scenes, where naive 3D coordinate clustering would fail, performing competitively with popular segmentation methods in unseen data in the single-view setting (see Sec. 4.3). We refer the reader to Appendix A.6.2 for implementation details and related visualizations.

4 Experiments

We design our experiments to explore the following questions: (i) **Sec. 4.1:** What are the contributions of our proposed object-centric priors for multi-view feature fusion? Does the dense number of views of our proposed dataset also contribute? (ii) **Sec. 4.2:** How does our method compare to state-of-the-art open-vocabulary approaches for semantic and referring segmentation tasks, both in multi- and in single-view settings? Is it robust to open-ended language? (iii) **Sec. 4.3:** What are the zero-shot generalization capabilities of our learned 3D representation in novel datasets that contain real-world scenes, as well as for the novel task of 3D instance segmentation? (v) **Sec. 4.4:** Can we leverage DROP-CLIP for language-guided 6-DoF robotic grasping?

4.1 Multi-view Feature Fusion Ablation Studies

To evaluate the contributions of our proposed object-centric priors, we conduct ablation studies on the multi-view feature fusion pipeline, where we compare 3D referring segmentation results of obtained 3D features in held-out scenes of MV-TOD. We highlight that here we aim to establish a performance **upper-bound** that the feature fusion method can provide for distillation, and not the distilled features themselves. We ablate: (i) patch-wise vs. object-wise fusion, (ii) MaskCLIP [8] patch-level vs. CLIP [7] masked crop

Fusion	f ^{2D}	$\Lambda_{v,i}$	$G_{v,i}$	mIoU	Ref.Segm (%)		
					Pr@25	Pr@50	Pr@75
point	patch	✓		37.3	55.4	33.7	16.7
point	patch		✓	57.0	74.1	59.5	40.9
point	patch	✓	✓	57.4	77.0	60.9	39.9
obj	obj			65.6	67.0	65.4	64.1
obj	obj	✓		67.3	68.7	67.1	65.8
obj	obj		✓	83.1	83.9	83.1	82.4
obj	obj	✓	✓	80.9	83.1	80.2	79.7

Table 2: Multi-view feature fusion ablation study for 3D referring segmentation in MV-TOD.

features, (iii) inclusion of visibility ($\Lambda_{v,i}$) and semantic informativeness ($G_{v,i}$) metrics for view selection. We report 3D segmentation metrics $mIoU$ and $Pr@X$ [39]. Results in Table 2.

Effect of object-centric priors We observe that all components contribute positively to the quality of the 3D features. Our proposed $G_{v,i}$ metric boosts $mIoU$ across both point- and object-wise fusion (57.0% vs. 44.2% and 83.1% vs. 65.6% respectively). Further, we observe that the usage of spatial priors for object-wise fusion and object-level features leads to drastic improvements, both in segmentation crispness (25.7% $mIoU$ delta), as well as in grounding precision (42.5% $Pr@75$ delta).

Effect of the number of views We ablate the 3D referring segmentation performance based on the number of input views in Fig. 5, where novel viewpoints are added incrementally.

We observe that in both setups (point- and object-wise) fusing features from more views leads to improvements, with a small plateauing behavior around 40 views. We believe this is an encouraging result for leveraging dense multi-view coverage in feature distillation pipelines, as we propose with MV-TOD. Please see Appendix A.3 for extended ablation studies that justify the design choices behind our fusion strategy, and Appendix A.5 for qualitative comparisons with vanilla approaches.

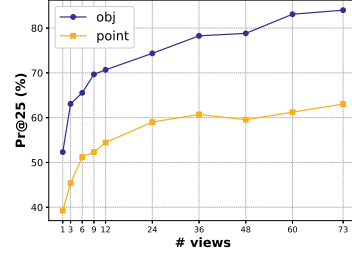


Figure 5: Referring segmentation precision vs. number of utilized views.

4.2 Open-Vocabulary 3D Segmentation Results in MV-TOD

In this section, we compare referring and semantic segmentation performance of our distilled features vs. previous open-vocabulary approaches, both in multi-view and in single-view settings.

For multi-view, we compare our trained model with OpenScene [12] and OpenMask3D [6] methods, where the full point-cloud from all 73 views is given as input. We note that for these baselines we obtain the upper-bound 3D features as before, as we observed that our trained model already outperforms them, so we refrained from also distilling features from baselines. For single-view, we feed our network with partial point-cloud from projected RGB-D pair, and compare with 2D baselines MaskCLIP [8] and OpenSeg [9] (see implementation details in Appendix A.4). Our model slightly outperforms the OpenMask3D upper bound baseline in the multi-view setting (+1.18% in referring and +2.57% in semantic segmentation), while significantly outperforming 2D baselines in the single-view setting (> 30% in both tasks). Importantly, single-view results closely match the multi-view ones ($\sim -4.0\%$), suggesting that DROP-CLIP indeed learns view-independent features. See Appendix A.5 for more qualitative comparisons with baselines.

Open-ended queries We evaluate the robustness of our model in different types of input language queries, organized in 4 families (class name - e.g. “cereal”, class + attribute - e.g. “brown cereal box”, open - e.g. “chocolate Kellogs”, and affordance - e.g. “I want something sweet”). Comparative results are presented in Fig. 6 and qualitative in Fig. 4. We observe that single-view performance closely follows that of upper-bound across query types, with multi-word affordance queries being the highest family of failures, potentially due to the “bag-of-words” behavior of CLIP text embeddings [16].

Method	#views	Ref.Segm. (%)				Sem.Segm (%)	
		mIoU	Pr@25	Pr@50	Pr@75	mIoU	mAcc ₂₅
OpenScene [†]	73	29.3	44.0	24.5	11.3	21.8	32.1
OpenMask3D ^{*†}	73	65.4	73.1	64.0	57.4	59.5	66.5
DROP-CLIP ^{*†}	73	82.7	86.1	82.4	79.2	75.4	80.0
DROP-CLIP	73	66.6	75.7	67.6	59.9	62.0	70.7
OpenSeg ^{→3D}	1	12.9	17.4	2.4	0.2	12.8	17.2
MaskCLIP ^{→3D}	1	25.6	40.4	18.7	7.0	21.0	32.1
DROP-CLIP	1	62.3	72.0	62.8	53.9	54.5	64.4

Table 3: Referring and Semantic segmentation results on MV-TOD test split. Methods with [†] denote upper-bound 3D features, whereas DROP-CLIP denotes our distilled model. Methods with ^{→3D} produce 2D predictions that are projected to 3D to compute metrics. Methods with * denote further usage of ground-truth segmentation masks.

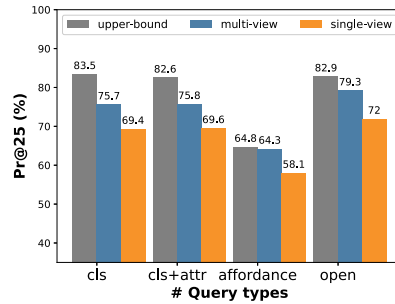


Figure 6: Referring segmentation precision vs. language query types.

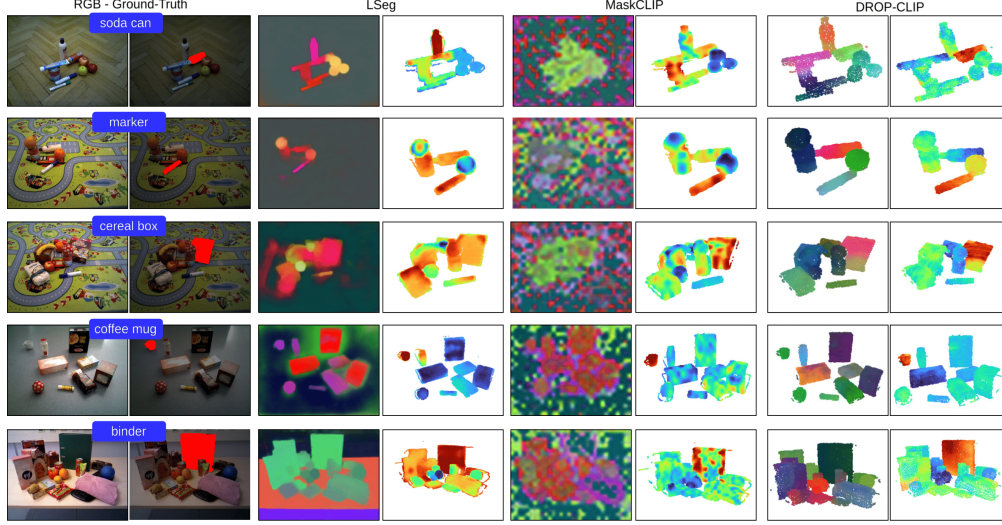


Figure 7: **Zero-Shot 3D Semantic Segmentation in Real Scenes:** Comparison of different referring segmentation models for five example cluttered indoor scenes from the OCID dataset. PCA features are displayed at pixel-level for 2D methods LSeg and MaskCLIP and in 3D for our point-cloud-based DROP-CLIP. Heatmaps from 2D models LSeg and MaskCLIP are projected to 3D for direct comparison with DROP-CLIP.

4.3 Generalization to Novel Domains / Tasks

Zero-shot transfer to real-world scenes In this section, we evaluate the zero-shot generalization capability of DROP-CLIP in real-world scenes that contain objects and vocabulary outside the MV-TOD distribution. We test in the validation split of the OCID-VLG [31] dataset, which contains 1249 queries from 165 unique cluttered tabletop scenes. We compare with 2D CLIP-based baselines LSeg [10], OpenSeg [9] and MaskCLIP [8] and popular 2D grounding method GroundedSAM [40] for the semantic segmentation task in the single-view setting as before.

Results are presented in Table 4. We find that even though fine-tuned in real data, baselines LSeg and OpenSeg under-perform compared to both MaskCLIP and our DROP-CLIP with a margin of $> 10\%$ mIoU, which we attribute to the distribution gap between the fine-tuning dataset ADE20K [41] and OCID scenes. These baselines tend to ground multiple regions in the scene, while MaskCLIP and DROP-CLIP provides tighter segmentations (see Fig. 7). When considering the stricter $mAcc_{75}$ metric, our approach scores a delta of 7.7% compared to MaskCLIP, suggesting a significant gain in grounding accuracy, especially in cases where the object is heavily occluded. Failures cases were observed in grounding objects that significantly vary in geometry and semantics from the MV-TOD catalog. Please see Appendix A.6 for further zero-shot experiments, comparisons with modern NeRF/3DGS methods and more qualitative results.

Zero-shot 3D instance segmentation We evaluate the potential of DROP-CLIP for out-of-the-box 3D instance segmentation via clustering the predicted features (see details in Appendix A.6.2). We conduct experiments for both the multi-view setting in MV-TOD, where we compare with Mask3D [42] transferred from the ScanRefer [1] checkpoint provided by the authors, where we feed full point-clouds from 73 views, as well as in OCID-VLG, where we compare with SAM [43] ViT-L model with single-view images. Results are summarized in Table 5. We observe that Mask3D struggles to

Method	OCID-VLG		
	<i>mIoU</i>	<i>mAcc</i> ₅₀	<i>mAcc</i> ₇₅
GroundedSAM	33.93	39.0	36.0
LSeg ^{→3D}	44.1	37.9	23.5
OpenSeg ^{→3D}	47.1	33.1	19.1
MaskCLIP ^{→3D}	57.1	59.4	31.0
DROP-CLIP	60.2	60.1	38.7

Table 4: Zero-shot semantic segmentation results (%) in the validation split of the OCID-VLG real-world dataset.

Method	OCID-VLG		MV-TOD	
	<i>mIoU</i>	<i>AP</i> ₂₅	<i>mIoU</i>	<i>AP</i> ₂₅
SAM	60.1	95.3	70.1	95.2
DROP-CLIP (S)	50.9	68.0	80.8	91.9
Mask3D	-	-	14.4	18.7
DROP-CLIP (F)	-	-	88.3	93.3

Table 5: Zero-shot 3D instance segmentation results in OCID-VLG (real-world) and our MV-TOD dataset.

generalize to tabletop domains, as it has been trained in room layout data with mostly furniture object categories. DROP-CLIP achieves an AP_{25} of 93.3%, illustrating that the learned 3D features can provide near-perfect instance segmentation in-distribution, even without explicit fine-tuning. When moving out-of-distribution in the single-view setting, we observe that DROP-CLIP achieves $mIoU$ that is competitive with foundation segmentation method SAM (50.9% vs. 60.1%). Failure cases include heavily cluttered regions of similar objects with same texture (e.g. food boxes), for which DROP-CLIP assigns very similar features that are identified as a single cluster.

4.4 Application: Language-guided Robotic Grasping

In this section, we wish to illustrate the applicability of DROP-CLIP in a language-guided robotic grasping scenario. We integrate our method with a 6-DoF grasp detection network [44], which proposes gripper poses for picking a target object segmented by DROP-CLIP. We randomly place 5-12 objects on a tabletop with different levels of clutter, and query the robot to pick a specific object, potentially amongst distractor objects of the same category. The user instruction is open-vocabulary and can involve open object descriptions, attributes, or user-affordances. We conducted 50 trials in Gazebo [45] and 10 with a real robot, and observed grounding accuracy of 84% and 80% respectively, and a final success rate of 64% and 60%. Motion failures were mostly due to grasp proposals for which the motion planning led to collisions. Similar to OCID, grounding failures were due to unseen query concepts and / or instances. Example trials are shown in Fig. 8, more details in Appendix A.7 and a robot demonstration video is provided as supplementary material.

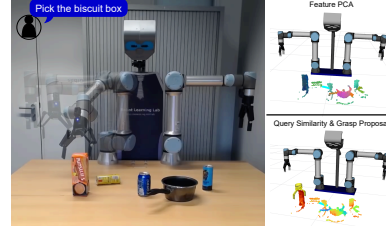


Figure 8: **Language-guided 6-DoF grasping:** Example robot trial (left), 3D features, grounding and grasp proposal (right).

5 Related work

We briefly discuss related efforts in this section, while a detailed comparison is given in Appendix A.8.

3D Scene Understanding There’s a long line of works in closed-set 3D scene understanding [46, 47, 48, 49, 50, 51], applied in 3D classification [52, 53], localization [54, 1] and segmentation [55, 23, 22], using two-stage pipelines with instance proposals from point-clouds [56, 57] or RGB-D views [58, 27], or single-stage methods [3] that leverage 3D-language cross attentions. [59] use CLIP embeddings for pretraining a 3D segmentation model, but still cannot be applied open-vocabulary.

Open-Vocabulary Grounding with CLIP Following the impressive results of CLIP [7] for open-set image recognition, followup works transfer CLIP’s powerful representations from image- to pixel-level [60, 61, 62, 63, 64, 65, 66, 9, 10, 8], extending to detection / segmentation, but limited to 2D. For 3D segmentation, the closest work is perhaps OpenMask3D [6] that extracts multi-view CLIP features from object proposals from Mask3D [42] to compute similarities with text queries.

3D CLIP Feature Distillation Recent works distill features from 2D foundation models with point-cloud encoders [12, 14, 21] or neural fields [13, 19, 17, 18, 19, 24], with applications in robot manipulation [20, 16] and navigation [67, 68]. However, associated works extract 2D features from OpenSeg [9], LSeg [10], MaskCLIP [8] or multi-scale crops from CLIP [7] and fuse point-wise with average pooling, while our approach leverages semantics-informed view selection and segmentation masks to do object-wise fusion with object-level features. Unlike all above field-based approaches, our method can be used real-time without the need for collecting multiple camera images at test-time.

6 Conclusion, Limitations & Future Work

We propose DROP-CLIP, a 2D→3D CLIP feature distillation framework that employs object-centric priors to select views based on semantic informativeness and ensure crisp 3D segmentations via

leveraging segmentation masks. Our method is designed to work from single-view RGB-D, encouraging view-independent features via distilling from dense multi-view scene coverage. We also release MV-TOD, a large-scale synthetic dataset of multi-view tabletop scenes with dense semantic / mask / grasp annotations. We believe our work can benefit the community, both in terms of released resources as well as illustrating and overcoming theoretical limitations of existing 3D feature distillation works.

While our spatial object-centric priors lead to improved segmentation quality, they collapse local features in favor of a global object-level feature, and hence cannot be applied for segmenting object parts. In the future, we plan to add object part annotations in our dataset and fuse with both object- and part-level masks. Second, DROP-CLIP cannot reconstruct 3D features that have significantly different geometry and / or semantics from the object catalog used during distillation. In the future we aim to explore modern generative text-to-3D models to further scale up the object and concept variety of MV-TOD. Finally, regarding robotic application, currently DROP-CLIP only provides language grounding, and a two-stage pipeline is necessary for robot grasping, while MV-TOD already provides rich 6-DoF grasp annotations. A next step would be to also distill them, opting for a joint 3D representation for grounding semantics and grasp affordances.

References

- [1] D. Z. Chen, A. X. Chang, and M. Nießner. Scanrefer: 3d object localization in rgb-d scans using natural language. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XX 16*, pages 202–221. Springer, 2020.
- [2] P. Achlioptas, A. Abdelreheem, F. Xia, M. Elhoseiny, and L. Guibas. Referit3d: Neural listeners for fine-grained 3d object identification in real-world scenes. *16th European Conference on Computer Vision (ECCV)*, 2020.
- [3] J. Luo, J. Fu, X. Kong, C. Gao, H. Ren, H. Shen, H. Xia, and S. Liu. 3d-sps: Single-stage 3d visual grounding via referred point progressive selection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16454–16463, 2022.
- [4] P.-H. Huang, H.-H. Lee, H.-T. Chen, and T.-L. Liu. Text-guided graph neural networks for referring 3d instance segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 1610–1618, 2021.
- [5] Z. Qian, Y. Ma, J. Ji, and X. Sun. X-refseg3d: Enhancing referring 3d instance segmentation via structured cross-modal graph neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 4551–4559, 2024.
- [6] A. Takmaz, E. Fedele, R. W. Sumner, M. Pollefeys, F. Tombari, and F. Engelmann. Openmask3d: Open-vocabulary 3d instance segmentation. *ArXiv*, abs/2306.13631, 2023. URL <https://api.semanticscholar.org/CorpusID:259243888>.
- [7] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from natural language supervision. *CoRR*, abs/2103.00020, 2021. URL <https://arxiv.org/abs/2103.00020>.
- [8] X. Dong, Y. Zheng, J. Bao, T. Zhang, D. Chen, H. Yang, M. Zeng, W. Zhang, L. Yuan, D. Chen, F. Wen, and N. Yu. Maskclip: Masked self-distillation advances contrastive language-image pretraining. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10995–11005, 2022. URL <https://api.semanticscholar.org/CorpusID:251799827>.
- [9] G. Ghiasi, X. Gu, Y. Cui, and T.-Y. Lin. Scaling open-vocabulary image segmentation with image-level labels. In *European Conference on Computer Vision*, 2021. URL <https://api.semanticscholar.org/CorpusID:250895808>.

- [10] B. Li, K. Q. Weinberger, S. J. Belongie, V. Koltun, and R. Ranftl. Language-driven semantic segmentation. *ArXiv*, abs/2201.03546, 2022. URL <https://api.semanticscholar.org/CorpusID:245836975>.
- [11] M. Oquab, T. Darcet, T. Moutakanni, H. Q. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza, F. Massa, A. El-Nouby, M. Assran, N. Ballas, W. Galuba, R. Howes, P.-Y. B. Huang, S.-W. Li, I. Misra, M. G. Rabbat, V. Sharma, G. Synnaeve, H. Xu, H. Jégou, J. Mairal, P. Labatut, A. Joulin, and P. Bojanowski. Dinov2: Learning robust visual features without supervision. *ArXiv*, abs/2304.07193, 2023. URL <https://api.semanticscholar.org/CorpusID:258170077>.
- [12] S. Peng, K. Genova, ChiyuMaxJiang, A. Tagliasacchi, M. Pollefeys, and T. A. Funkhouser. Openscene: 3d scene understanding with open vocabularies. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 815–824, 2022. URL <https://api.semanticscholar.org/CorpusID:254044069>.
- [13] J. Kerr, C. M. Kim, K. Goldberg, A. Kanazawa, and M. Tancik. Lerf: Language embedded radiance fields. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 19672–19682, 2023. URL <https://api.semanticscholar.org/CorpusID:257557329>.
- [14] P. D. Nguyen, T. Ngo, C. Gan, E. Kalogerakis, A. D. Tran, C. Pham, and K. Nguyen. Open3dis: Open-vocabulary 3d instance segmentation with 2d mask guidance. *ArXiv*, abs/2312.10671, 2023. URL <https://api.semanticscholar.org/CorpusID:266348609>.
- [15] S. Koch, N. Vaskevicius, M. Colosi, P. Hermosilla, and T. Ropinski. Open3dsg: Open-vocabulary 3d scene graphs from point clouds with queryable objects and open-set relationships. *ArXiv*, abs/2402.12259, 2024. URL <https://api.semanticscholar.org/CorpusID:267750890>.
- [16] B. W. Shen, G. Yang, A. Yu, J. R. Wong, L. P. Kaelbling, and P. Isola. Distilled feature fields enable few-shot language-guided manipulation. In *Conference on Robot Learning*, 2023. URL <https://api.semanticscholar.org/CorpusID:260926035>.
- [17] V. Tschernezki, I. Laina, D. Larlus, and A. Vedaldi. Neural feature fusion fields: 3d distillation of self-supervised 2d image representations. *2022 International Conference on 3D Vision (3DV)*, pages 443–453, 2022. URL <https://api.semanticscholar.org/CorpusID:252118532>.
- [18] S. Kobayashi, E. Matsumoto, and V. Sitzmann. Decomposing nerf for editing via feature field distillation. *ArXiv*, abs/2205.15585, 2022. URL <https://api.semanticscholar.org/CorpusID:249209811>.
- [19] F. Engelmann, F. Manhardt, M. Niemeyer, K. Tateno, M. Pollefeys, and F. Tombari. Opennerf: Open set 3d neural scene segmentation with pixel-wise features and rendered novel views, 2024.
- [20] A. Rashid, S. Sharma, C. M. Kim, J. Kerr, L. Y. Chen, A. Kanazawa, and K. Goldberg. Language embedded radiance fields for zero-shot task-oriented grasping. In *Conference on Robot Learning*, 2023. URL <https://api.semanticscholar.org/CorpusID:261882332>.
- [21] J. Zhang, R. Dong, and K. Ma. Clip-fo3d: Learning free open-world 3d scene representations from 2d dense clip. *2023 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, pages 2040–2051, 2023. URL <https://api.semanticscholar.org/CorpusID:257404908>.
- [22] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5828–5839, 2017.

- [23] S. K. Ramakrishnan, A. Gokaslan, E. Wijmans, O. Maksymets, A. Clegg, J. Turner, E. Undersander, W. Galuba, A. Westbury, A. X. Chang, M. Savva, Y. Zhao, and D. Batra. Habitat-matterport 3d dataset (hm3d): 1000 large-scale 3d environments for embodied ai. *ArXiv*, abs/2109.08238, 2021. URL <https://api.semanticscholar.org/CorpusID:237563216>.
- [24] M. Qin, W. Li, J. Zhou, H. Wang, and H. Pfister. Langsplat: 3d language gaussian splatting, 2024.
- [25] M. Chen, Q. Hu, Z. Yu, H. Thomas, A. Feng, Y. Hou, K. McCullough, F. Ren, and L. Soibelman. Stpls3d: A large-scale synthetic and real aerial photogrammetry 3d point cloud dataset. *arXiv preprint arXiv:2203.09065*, 2022.
- [26] J. Straub, T. Whelan, L. Ma, Y. Chen, E. Wijmans, S. Green, J. J. Engel, R. Mur-Artal, C. Ren, S. Verma, et al. The replica dataset: A digital replica of indoor spaces. *arXiv preprint arXiv:1906.05797*, 2019.
- [27] H. Liu, A. Lin, X. Han, L. Yang, Y. Yu, and S. Cui. Refer-it-in-rgbd: A bottom-up approach for 3d visual grounding in rgbd images. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6028–6037, 2021.
- [28] C. Mauceri, M. Palmer, and C. Heckman. Sun-spot: An rgb-d dataset with spatial referring expressions. *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*, pages 1883–1886, 2019.
- [29] H.-S. Fang, C. Wang, M. Gou, and C. Lu. Graspnet-1billion: A large-scale benchmark for general object grasping. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11444–11453, 2020.
- [30] H. Zhang, D. Yang, H. Wang, B. Zhao, X. Lan, J. Ding, and N. Zheng. Regrad: A large-scale relational grasp dataset for safe and object-specific robotic grasping in clutter. *IEEE Robotics and Automation Letters*, 7(2):2929–2936, 2022.
- [31] G. Tzifas, X. Yucheng, A. Goel, M. Kasaei, Z. Li, and H. Kasaei. Language-guided robot grasping: Clip-based referring grasp synthesis in clutter. In *7th Annual Conference on Robot Learning*, 2023.
- [32] A. D. Vuong, M. N. Vu, H. Le, B. Huang, B. P. K. Huynh, T. D. Vo, A. Kugi, and A. Nguyen. Grasp-anything: Large-scale grasp dataset from foundation models. *ArXiv*, abs/2309.09818, 2023. URL <https://api.semanticscholar.org/CorpusID:262045996>.
- [33] I. Armeni, O. Sener, A. R. Zamir, H. Jiang, I. Brilakis, M. Fischer, and S. Savarese. 3d semantic parsing of large-scale indoor spaces. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1534–1543, 2016.
- [34] D. Rozenberszki, O. Litany, and A. Dai. Language-grounded indoor 3d semantic segmentation in the wild. In *European Conference on Computer Vision*, pages 125–141. Springer, 2022.
- [35] C. Eppner, A. Mousavian, and D. Fox. ACRONYM: A large-scale grasp dataset based on simulation. In *2021 IEEE Int. Conf. on Robotics and Automation, ICRA*, 2020.
- [36] B. O. Community. Blender - a 3d modelling and rendering package. 2018. URL <http://www.blender.org>.
- [37] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, J. Xiao, L. Yi, and F. Yu. ShapeNet: An Information-Rich 3D Model Repository. Technical Report arXiv:1512.03012 [cs.GR], Stanford University — Princeton University — Toyota Technological Institute at Chicago, 2015.

- [38] Gpt-4v(ision) system card. 2023. URL <https://api.semanticscholar.org/CorpusID:263218031>.
- [39] C. Wu, Y. Liu, J. Ji, Y. Ma, H. Wang, G. Luo, H. Ding, X. Sun, and R. Ji. 3d-gres: Generalized 3d referring expression segmentation. *ArXiv*, abs/2407.20664, 2024. URL <https://api.semanticscholar.org/CorpusID:271544474>.
- [40] T. Ren, S. Liu, A. Zeng, J. Lin, K. Li, H. Cao, J. Chen, X. Huang, Y. Chen, F. Yan, Z. Zeng, H. Zhang, F. Li, J. Yang, H. Li, Q. Jiang, and L. Zhang. Grounded sam: Assembling open-world models for diverse visual tasks, 2024.
- [41] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba. Scene parsing through ade20k dataset. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5122–5130, 2017. URL <https://api.semanticscholar.org/CorpusID:5636055>.
- [42] J. Schult, F. Engelmann, A. Hermans, O. Litany, S. Tang, and B. Leibe. Mask3D: Mask Transformer for 3D Semantic Instance Segmentation. 2023.
- [43] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, P. Dollár, and R. B. Girshick. Segment anything. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3992–4003, 2023. URL <https://api.semanticscholar.org/CorpusID:257952310>.
- [44] S. Chen, W. N. Tang, P. Xie, W. Yang, and G. Wang. Efficient heatmap-guided 6-dof grasp detection in cluttered scenes. *IEEE Robotics and Automation Letters*, 8:4895–4902, 2023. URL <https://api.semanticscholar.org/CorpusID:259363869>.
- [45] N. P. Koenig and A. Howard. Design and use paradigms for gazebo, an open-source multi-robot simulator. *2004 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) (IEEE Cat. No.04CH37566)*, 3:2149–2154 vol.3, 2004.
- [46] C. B. Choy, J. Gwak, and S. Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3070–3079, 2019. URL <https://api.semanticscholar.org/CorpusID:121123422>.
- [47] L. Han, T. Zheng, L. Xu, and L. Fang. Occuseg: Occupancy-aware 3d instance segmentation. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2937–2946, 2020. URL <https://api.semanticscholar.org/CorpusID:212725768>.
- [48] W. Hu, H. Zhao, L. Jiang, J. Jia, and T.-T. Wong. Bidirectional projection network for cross dimension scene understanding. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14368–14377, 2021. URL <https://api.semanticscholar.org/CorpusID:232379958>.
- [49] Z. Hu, X. Bai, J. Shang, R. Zhang, J. Dong, X. Wang, G. Sun, H. Fu, and C.-L. Tai. Vmnet: Voxel-mesh network for geodesic-aware 3d semantic segmentation. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 15468–15478, 2021. URL <https://api.semanticscholar.org/CorpusID:236493200>.
- [50] J. Li, X. He, Y. Wen, Y. Gao, X. Cheng, and D. Zhang. Panoptic-phnet: Towards real-time and high-precision lidar panoptic segmentation via clustering pseudo heatmap. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11799–11808, 2022. URL <https://api.semanticscholar.org/CorpusID:248811224>.
- [51] D. Robert, B. Vallet, and L. Landrieu. Learning multi-view aggregation in the wild for large-scale 3d semantic segmentation. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5565–5574, 2022. URL <https://api.semanticscholar.org/CorpusID:248218804>.

- [52] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao. 3d shapenets: A deep representation for volumetric shapes. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1912–1920, 2014. URL <https://api.semanticscholar.org/CorpusID:206592833>.
- [53] R. Zhang, Z. Guo, W. Zhang, K. Li, X. Miao, B. Cui, Y. J. Qiao, P. Gao, and H. Li. Pointclip: Point cloud understanding by clip. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8542–8552, 2021. URL <https://api.semanticscholar.org/CorpusID:244909021>.
- [54] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom. nuscnets: A multimodal dataset for autonomous driving. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11618–11628, 2019. URL <https://api.semanticscholar.org/CorpusID:85517967>.
- [55] J. Behley, M. Garbade, A. Milioto, J. Quenzel, S. Behnke, C. Stachniss, and J. Gall. SemanticKITTI: A dataset for semantic scene understanding of lidar sequences. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9296–9306, 2019. URL <https://api.semanticscholar.org/CorpusID:199441943>.
- [56] P. Achlioptas, A. Abdelreheem, F. Xia, M. Elhoseiny, and L. J. Guibas. Referit3d: Neural listeners for fine-grained 3d object identification in real-world scenes. In *European Conference on Computer Vision*, 2020.
- [57] L. Zhao, D. Cai, L. Sheng, and D. Xu. 3dvg-transformer: Relation modeling for visual grounding on point clouds. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2908–2917, 2021.
- [58] S. Huang, Y. Chen, J. Jia, and L. Wang. Multi-view transformer for 3d visual grounding. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15503–15512, 2022.
- [59] D. Rozenberszki, O. Litany, and A. Dai. Language-grounded indoor 3d semantic segmentation in the wild. *ArXiv*, abs/2204.07761, 2022. URL <https://api.semanticscholar.org/CorpusID:248227627>.
- [60] X. Gu, T.-Y. Lin, W. Kuo, and Y. Cui. Open-vocabulary object detection via vision and language knowledge distillation. In *International Conference on Learning Representations*, 2021. URL <https://api.semanticscholar.org/CorpusID:238744187>.
- [61] Y. Zhong, J. Yang, P. Zhang, C. Li, N. C. F. Codella, L. H. Li, L. Zhou, X. Dai, L. Yuan, Y. Li, and J. Gao. Regionclip: Region-based language-image pretraining. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16772–16782, 2021. URL <https://api.semanticscholar.org/CorpusID:245218534>.
- [62] M. Minderer, A. A. Gritsenko, A. Stone, M. Neumann, D. Weissenborn, A. Dosovitskiy, A. Mahendran, A. Arnab, M. Dehghani, Z. Shen, X. Wang, X. Zhai, T. Kipf, and N. Houlsby. Simple open-vocabulary object detection with vision transformers. *ArXiv*, abs/2205.06230, 2022. URL <https://api.semanticscholar.org/CorpusID:248721818>.
- [63] X. Zhou, R. Girdhar, A. Joulin, P. Krahenbuhl, and I. Misra. Detecting twenty-thousand classes using image-level supervision. *ArXiv*, abs/2201.02605, 2022. URL <https://api.semanticscholar.org/CorpusID:245827815>.
- [64] M. Minderer, A. A. Gritsenko, and N. Houlsby. Scaling open-vocabulary object detection. *ArXiv*, abs/2306.09683, 2023. URL <https://api.semanticscholar.org/CorpusID:259187664>.

- [65] Z. Wang, Y. Lu, Q. Li, X. Tao, Y. Guo, M. Gong, and T. Liu. Cris: Clip-driven referring image segmentation. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11676–11685, 2021. URL <https://api.semanticscholar.org/CorpusID:244729320>.
- [66] T. Lüddecke and A. S. Ecker. Image segmentation using text and image prompts. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7076–7086, 2021. URL <https://api.semanticscholar.org/CorpusID:247794227>.
- [67] N. M. M. Shafiullah, C. Paxton, L. Pinto, S. Chintala, and A. Szlam. Clip-fields: Weakly supervised semantic fields for robotic memory. *ArXiv*, abs/2210.05663, 2022. URL <https://api.semanticscholar.org/CorpusID:252815898>.
- [68] B. Bolte, A. S. Wang, J. Yang, M. Mukadam, M. Kalakrishnan, and C. Paxton. Usa-net: Unified semantic and affordance representations for robot memory. *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1–8, 2023. URL <https://api.semanticscholar.org/CorpusID:258298248>.
- [69] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. 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.
- [70] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan. Supervised contrastive learning. *CoRR*, abs/2004.11362, 2020. URL <https://arxiv.org/abs/2004.11362>.
- [71] H. Oki, M. Abe, J. Miyao, and T. Kurita. Triplet loss for knowledge distillation. *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–7, 2020. URL <https://api.semanticscholar.org/CorpusID:215814195>.
- [72] L. Yang, Y. Wang, X. Li, X. Wang, and J. Yang. Fine-grained visual prompting. *ArXiv*, abs/2306.04356, 2023. URL <https://api.semanticscholar.org/CorpusID:259096008>.
- [73] A. Shtedritski, C. Rupprecht, and A. Vedaldi. What does clip know about a red circle? visual prompt engineering for vlms. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11953–11963, 2023. URL <https://api.semanticscholar.org/CorpusID:258108138>.
- [74] S. Ainetter and F. Fraundorfer. End-to-end trainable deep neural network for robotic grasp detection and semantic segmentation from rgb. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages 13452–13458, 2021.
- [75] J. Guo, X. Ma, Y. Fan, H. Liu, and Q. Li. Semantic gaussians: Open-vocabulary scene understanding with 3d gaussian splatting. *ArXiv*, abs/2403.15624, 2024. URL <https://api.semanticscholar.org/CorpusID:268680548>.
- [76] R.-Z. Qiu, G. Yang, W. Zeng, and X. Wang. Feature splatting: Language-driven physics-based scene synthesis and editing. *ArXiv*, abs/2404.01223, 2024. URL <https://api.semanticscholar.org/CorpusID:268819312>.
- [77] S. Zhou, H. Chang, S. Jiang, Z. Fan, Z. Zhu, D. Xu, P. Chari, S. You, Z. Wang, and A. Kadambi. Feature 3dgs: Supercharging 3d gaussian splatting to enable distilled feature fields. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 21676–21685, 2023. URL <https://api.semanticscholar.org/CorpusID:265722936>.

- [78] R. Ding, J. Yang, C. Xue, W. Zhang, S. Bai, and X. Qi. Pla: Language-driven open-vocabulary 3d scene understanding. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7010–7019, 2022. URL <https://api.semanticscholar.org/CorpusID:254069374>.
- [79] J. Yang, R. Ding, Z. Wang, and X. Qi. Regionplc: Regional point-language contrastive learning for open-world 3d scene understanding. *ArXiv*, abs/2304.00962, 2023. URL <https://api.semanticscholar.org/CorpusID:257913360>.
- [80] R. Ding, J. Yang, C. Xue, W. Zhang, S. Bai, and X. Qi. Lowis3d: Language-driven open-world instance-level 3d scene understanding. *IEEE transactions on pattern analysis and machine intelligence*, PP, 2023. URL <https://api.semanticscholar.org/CorpusID:260351247>.
- [81] D. Hegde, J. M. J. Valanarasu, and V. M. Patel. Clip goes 3d: Leveraging prompt tuning for language grounded 3d recognition. *2023 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, pages 2020–2030, 2023. URL <https://api.semanticscholar.org/CorpusID:257632366>.
- [82] Z. Huang, X. Wu, X. Chen, H. Zhao, L. Zhu, and J. Lasenby. Openins3d: Snap and lookup for 3d open-vocabulary instance segmentation. *ArXiv*, abs/2309.00616, 2023. URL <https://api.semanticscholar.org/CorpusID:261494064>.
- [83] S. Lu, H. Chang, E. P. Jing, A. Boularias, and K. E. Bekris. Ovir-3d: Open-vocabulary 3d instance retrieval without training on 3d data. *ArXiv*, abs/2311.02873, 2023. URL <https://api.semanticscholar.org/CorpusID:262072783>.
- [84] Y. nuo Yang, X. Wu, T. He, H. Zhao, and X. Liu. Sam3d: Segment anything in 3d scenes. *ArXiv*, abs/2306.03908, 2023. URL <https://api.semanticscholar.org/CorpusID:259088699>.
- [85] M. Yan, J. Zhang, Y. Zhu, and H. R. Wang. Maskclustering: View consensus based mask graph clustering for open-vocabulary 3d instance segmentation. *ArXiv*, abs/2401.07745, 2024. URL <https://api.semanticscholar.org/CorpusID:266999755>.
- [86] Y. Yin, Y. Liu, Y. Xiao, D. Cohen-Or, J. Huang, and B. Chen. Sai3d: Segment any instance in 3d scenes. *ArXiv*, abs/2312.11557, 2023. URL <https://api.semanticscholar.org/CorpusID:266362709>.
- [87] W. Dai, J. Li, D. Li, A. M. H. Tiong, J. Zhao, W. Wang, B. A. Li, P. Fung, and S. C. H. Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *ArXiv*, abs/2305.06500, 2023. URL <https://api.semanticscholar.org/CorpusID:258615266>.
- [88] B. Kerbl, G. Kopanas, T. Leimkuehler, and G. Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (TOG)*, 42:1 – 14, 2023. URL <https://api.semanticscholar.org/CorpusID:259267917>.