

A Appendix

A.1 MV-TOD Details

In this section we provide details for generating our MV-TOD scenes and their annotations (Sec. A.1.1) and present some statistics for the object and query catalog of MV-TOD (Sec. A.1.2).

A.1.1 Dataset Generation

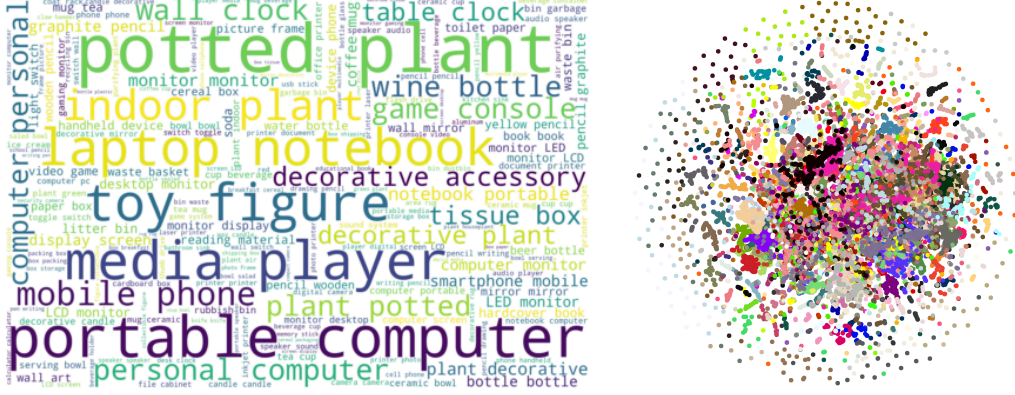


Figure 9: A wordcloud and T-SNE embedding projection visualization of textual concepts included in MV-TOD.

We generate the MV-TOD dataset in Blender [36] engine with following steps:

Random object spawn For each scene, firstly, a support plane is spawned at the origin position. Then, random objects are selected to set up the multi-object tabletop scene. The number of objects ranges from 4 to 12, to make sure that our dataset covers both isolated and cluttered scenes. All selected objects are then spawned above the support plane with random position and random rotation. It is important to note that due to the limitation of Blender physical engine, an additional collision check is needed when an object is spawned into the scene to avoid initial collision. Since Blender does not provide users with the APIs to do the collision check, we check the collision by calculating the 3D IoU between object bounding boxes. After spawning object, the internal physical simulator is launched to simulation the falling of all the spawned object onto the plane. Once the objects are still, the engine will start rendering images.

Multi-view rendering In total 73 cameras are set in each scene for rendering images from different views. One of them are spawned right on top of the origin position for rendering a top-down image, while the rest are uniformly distributed on the surface of the top hemisphere. An RGB image, a depth image (with the raw depth information in meters), and an instance segmentation mask are rendered at each view. All the annotations are saved in the COCO [69] JSON format for each scene.

Data augmentation In order to diversify the generated data, several augmentation methods are applied. Firstly, different textures and materials are randomly applied to the support plane, as well as the scene background, to simulate different types of table surfaces and background environments. Second, when the objects are spawned, their sizes and materials are randomly jittered. Thirdly, we also randomly slightly modify the position of cameras towards the radial direction. Finally, the position and intensity of the light object in each scene are also randomly set.

Semantic object annotation generation To offer the functionality of querying target objects in our dataset by using high-level concepts and distinguish similar objects using fine-grained attributes, we also provide per-object semantic concepts generated with the aid of large vision-language models. For each object CAD model, we render 10 observation images from different views in Blender. Then, these images, together with an instruction prompt are fed to GPT-4V [38] to generate a response

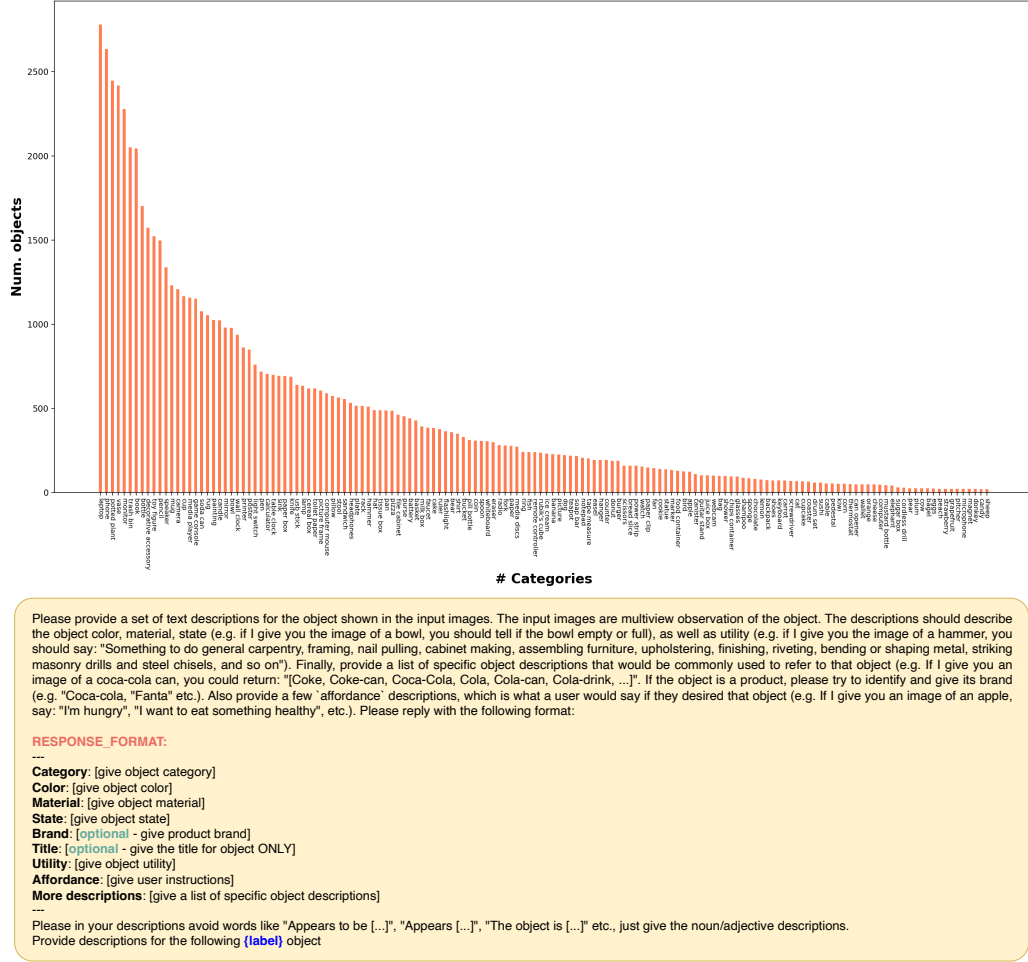


Figure 11: The text prompt we used for instructing GPT-4V. The **{label}** token will be replaced by the class name of the current object.

describing the current object in different perspectives, including category, color, material, state, utility, affordance, title (if applicable), and brand (if applicable). The text prompt we used to instruct GPT-4V is presented in Figure 11.

6-DoF grasp annotations Since our model set originates from ShapeNet-Sem [37], we leverage the object-wise 6-DoF grasp proposals generated previously in the ACRONYM dataset [35]. These grasps were executed and evaluated in a simulation environment, leading to a total of 2000 grasp candidates per object. We filter the successful grasps and connect them with each object instance in each of our scenes, by transforming the grasp annotation according to the recorded object’s 6D pose from Blender. We further filter grasps by rendering a gripper mesh and removing all grasp poses that lead to collisions with the table or other objects.

A.1.2 Dataset Analysis

We visualize a wordcloud of the concept vocabulary of MV-TOD, together with tSNE projections of their CLIP text embeddings in Figure 9. Certain object names (e.g. "plant", "computer", "phone", "vase") appear more frequently, as those are the objects that are most frequent in MV-TOD object catalog, hence they spawn a lot of expressions referring to them. Besides common class names, the wordcloud demonstrates that the most frequent concepts used to disambiguate objects are

Type	Train	Test
Class	66.8k	19.2k
Class+Attr	76.5k	21.8k
Affordance	151.1k	44.7k
Open	356.8k	102.1k

Table 6: Number of referring expressions in MV-TOD organized by type

supplementary attributes (e.g. *decorative*, *potted*, *portable*, etc). Finally, colors and materials appear also frequently, as they are a common discriminating attribute between objects of the same category.

We further provide statistical analysis of MV-TOD in Table 6 and Figure 10. The number of referring expressions categorized by their types are listed in Table 6. We provide rich *open* expressions, which stems from open vocabulary concepts that can describe the referred objects in various aspects. As it can be seen Figure 10, there exists a typical long-tail distribution in our dataset in terms of the number of objects per-category, where *laptop*, *phone*, and *plant* have the most variant instances.

A.2 Distillation Implementation Details

We use the ViT-L/14@336px variant of CLIP’s vision encoder, which provides features of size $C = 768$ from 336×448 image inputs with patch size 14. We distill with a MinkowskiNet14D [46] sparse 3D-UNet backbone, which consists of 8 sparse ResNet blocks with output sizes of (32, 64, 128, 256, 384, 384, 384, 384) and a final 1×1 convolution head to 768 channels. To increase the 3D coordinates resolution, we upscale the input point-clouds to $\times 10$ and voxelize with original dimension of $d = 0.02$ (for feature fusion), and a voxel grid $d = 0.05$ for training with the Minkowski framework. To reduce the input dimensionality and speedup training and inference time, we remove the table points via filtering out the table’s 3D mask.¹ We train using AdamW with initial learning rate $3 \cdot 10^{-4}$ and cosine annealing to 10^{-4} over 300 epochs, and a weight decay of 10^{-4} . In each ResNet block, we include sparse batch normalization layers with momentum of 0.1. We train using two RTX 4090 GPUs, which takes about 4 days. Following [12, 6] we use spatial augmentations such as elastic distortion, horizontal flipping and small random translations and rotations. We also employ color-based augmentation such as chromatic auto-contrast, random color translation, jitter and hue saturation translation. To better emulate partial views with greater diversity, we train with full point-clouds but further add a per-object blob removal augmentation method that removes consistent blobs of points from each object instance. After training for 300 epochs, we fine-tune our obtained checkpoint on only partial point-clouds from randomly sampled views for each scene in our dataset. We experimented with several auxiliary losses to reinforce within-object feature similarity, such as supervised contrastive loss [70] as well as KL triplet loss [71],

Hyper-parameter	Value
voxel_size	0.05
feat_dim	768
color_trans_ratio	0.01
color_jitter_std	0.02
hue_max	0.01
saturation_max	0.1
elastic_distortion_granularity_min	0.1
elastic_distortion_granularity_max	0.3
elastic_distortion_magnitude_min	0.4
elastic_distortion_magnitude_max	0.8
n_blob_min	1
n_blob_max	2
blob_size_min	50
blob_size_max	101
random_euler_order	True
random_rot_chance	0.6
rotate_min_x	-0.1309
rotate_max_x	0.1309
rotate_min_y	-0.1309
rotate_max_y	0.1309
rotate_min_z	-0.1309
rotate_max_z	0.1309
arch_3d	MinkUNet14D
batch_size	8
batch_size_val	8
base_lr	0.0003
weight_decay	0.00001
min_lr	0.0001
loss_type	cosine
use_aux_loss	False
use_cls_head	False
loss_weight_aux	1.0
loss_weight_cls	0.1
dropout_rate	0.0
epochs	300
power	0.9
momentum	0.9
max_norm	5.0
sync_bn	True

Table 7: Training hyper-parameters

¹. During inference, we employ RANSAC to remove table points without access to segmentation masks.

but found that they do not significantly contribute to convergence compared to using only the main cosine distance loss. See Table 7 for a full overview of training and augmentation hyper-parameters.

A.3 Extended Multi-View Feature Fusion Ablations

Our object-centric fusion pipeline considers several design choices besides the ones discussed in the main paper. In particular, we study: (i) Why masked crops as input to object-level 2D CLIP feature computation? How does it compare with other popular visual prompts to CLIP?, (ii) Why equations (3) and (4) in the semantic informativeness metric computation? How to sample negatives?, and (iii) What is the best strategy and hyper-parameters for doing inference?

CLIP visual prompts Previous works have extensively studied how to prompt CLIP to make it focus in a particular entity in the scene [72, 73]. We study the potential of visual prompting for obtaining object-level CLIP features in our object-centric feature fusion pipeline, via measuring their final referring segmentation *mIoU* in a subset of MV-TOD validation split. We compile the following visual prompt options: (a) *crop*, where we crop a bounding box around each

crop	crop-mask	mask-blur	mask-gray	mask-out	#crops	crop-ratio	mIoU (%)
×					1	-	81.9
×					3	0.1	81.0
×					3	0.15	80.6
×					3	0.2	80.3
	×				1	-	84.0
	×				3	0.15	82.4
	×				3	0.2	82.0
×	×				1	-	81.7
		×			-	-	74.6
			×		-	-	57.4
				×	-	-	79.7
		×	×	×	-	-	70.2

Table 8: CLIP visual prompt ablation studies.

object [13, 6], (b) *crop-mask*, where we crop a bounding box but only leave the pixels of the object’s 2D instance mask inside and uniform paint the background (black, white or gray, based on the mask’s mean color), (c) *mask-{blur, gray, out}*, where we use the entire image with the target object instance highlighted [72] and the rest completely removed as before (*out*), converted to grayscale (*gray*) or applied a median blur filter (*blur*). For the crop options, we further ablate the number of multi-scale crops used and their relative expansion ratio. Results are summarized in Table 8. We observe the following: (a) image-level visual prompts used previously [72] do not perform as well as cropped bounding boxes, (b) using multi-scale crops [13, 6] doesn’t improve over using a single object crop, (c) masked crops outperform non-masked crops by a small margin of 2.1%. The difference is due to cases of heavy clutter, when the bounding box of the non-masked crop also includes neighboring objects, making the representation obtained by CLIP also give high similarities with the neighbor’s prompt. This effect is more pronounced when using multiple crops with larger expansion ratios, as more and more neighboring objects are included in the crops.

Semantic informativeness metric We ablate the following components when computing semantic informativeness metric $G_{v,n}$: (i) the type of prompts used as q^+, Q^- , i.e. *cls* for category-level prompts and *open* where we use all instance-level descriptions annotated with GPT-4V, (ii) the operator used to reduce the negative prompts to single feature dimension, i.e. *max* and *mean*, and (iii) how to sample negatives for Q^- , i.e. including only negative prompts for objects in the *scene*, or including *all* other dataset objects. Results are shown in Table 9. First, we observe that *max* operator generally outperforms *mean*, with the exception of when using *all* negatives. However, the best configuration was using *max* operator with *scene* negatives. Second, using *open* prompts provides marginal improvements over *cls* in all other settings. Finally, using *scene* negatives outperforms using *all* in most cases. This is because when using all negatives from the dataset, some semantic concepts will be highly similar with the positive prompt, making the metric too ‘strict’, as only few views will pass the condition $G_{v,n} \geq 0$.

Prompts	Operator	Negatives	Ref.Segm. (%)		Sem.Segm (%)	
			mIoU	Pr@25	mIoU	mAcc
cls	mean	scene	82.2	83.1	73.1	75.1
cls	max	scene	82.8	84.0	74.9	76.7
cls	mean	all	80.9	81.0	71.6	73.9
cls	max	all	72.6	75.1	60.7	63.2
open	mean	scene	76.4	78.8	68.3	70.3
open	max	scene	83.9	85.5	75.6	77.2
open	mean	all	81.0	81.8	71.6	74.0
open	max	all	72.9	74.2	63.8	64.6

Table 9: Semantic informativeness metric ablation studies. Results in MV-TOD validation subset.

Inference strategies As discussed in Sec. 3.4 there are two methods for performing referring segmentation inference: (a) selecting all points with higher probability for positive vs. maximum negative prompt $\rho^+ > \mathcal{P}^-$, or (b) thresholding ρ^+ with a hyper-parameter s_{thr} . Additionally, we compare the final referring segmentation performance based on the negative prompts used at test time: (a) prompts from object instances within the *scene*, (b) prompts from *all* dataset object instances (similar to semantic segmentation task), (c) fixed *canonical* phrases

{“*object*”, “*thing*”, “*texture*”, “*stuff*”} [13], and (d) no negative prompts (-), where we threshold the raw cosine similarities with the positive query. Results in Table 10. We observe that thresholding provides better results than the first method when the right threshold is chosen, a result which we found holds also for our distilled model. A high threshold of 0.95 was found optimal for upper bound experiments, while a threshold of 0.7 for our distilled model, although we further fine-tuned it for zero-shot and robot experiments (see Sec. A.6). Regarding negative prompts, as expected, providing in-scene negatives gives the best results, with a significant delta from canonical (7.6%), all (7.9%) and no negatives (12.6%). However, we observe that even without such prior, the performance is still competitive, even when entirely skipping negative prompts.

Method	Negatives	Ref.Segm. (%)			
		mIoU	Pr@25	Pr@50	Pr@75
$\rho^+ > \mathcal{P}^-$	scene	73.7	77.4	73.0	69.8
$\rho^+ > \mathcal{P}^-$	canonical	53.4	57.4	52.6	49.7
$\rho^+ > \mathcal{P}^-$	all	30.8	31.0	31.0	30.8
$s_{thr}@0.95$	scene	82.8	84.0	83.2	82.0
$s_{thr}@0.95$	canonical	75.2	77.6	74.7	72.9
$s_{thr}@0.95$	all	74.9	76.6	75.4	73.0
$s_{thr}@0.95$	-	70.2	70.6	69.9	69.5
$s_{thr}@0.9$	scene	82.1	83.6	82.8	79.8
$s_{thr}@0.8$	scene	79.9	83.0	80.4	75.7

Table 10: Inference method ablation studies.

A.4 Baseline Implementations

OpenSeg [9] extends CLIP’s image-level visual representations to pixel-level, by first proposing instance segmentation masks and then aligning them to matched text captions. Given a text query, with OpenSeg we can obtain a 2D instance segmentation mask. For extending to 3D, we project the 2D mask pixels to 3D according to the mask region’s depth values and camera intrinsics and transform to world frame.

LSeg [10] similarly trains an image encoder to be aligned with CLIP text embeddings at pixel-level with dense contranstive loss, therefore allowing open-vocabulary queries at test-time. Similar to OpenSeg, we project 2D predictions to 3D according to depth and camera intrinsics and transform to world frame to compute metrics.

MaskCLIP [8] provides a drop-in reparameterization trick in the attention pooling layer of CLIP’s ViT encoder, enabling text-aligned patch features that can be directly used for grounding tasks. We use bicubic interpolation to upsample the patch-level features to pixel-level before computing cosine similarities with text queries. Similar to OpenSeg and LSeg, we project and transform the predicted 2D mask to calculate 3D metrics.

OpenScene [12] is the first method to introduce the 3D feature distillation methodology for room scan datasets. It utilizes OpenSeg [9] to extract pixel-level 2D features and fuses them point-wise with vanilla average pooling, as formulated in Sec. 3.1. To provide fair comparisons with our approach, and as we found that MaskCLIP’s features perform favourably vs. OpenSeg’s, we use patch-wise MaskCLIP features, interpolated to original image size. We aggregate all 73 views, perform vanilla feature fusion in the full point-cloud, and measure the final fused 3D feature’s performance as the OpenScene performance. We highlight that this setup represents the *upper-bound* performance OpenScene can provide, as we use the target 3D features and not distilled ones obtained through training, which we refrained from doing, as our results already outperform OpenScene’s upper bound.

OpenMask3D [6] is a recent two-stage method for referring segmentation in point-cloud data. In the first stage, Mask3D [42] is used for 3D instance segmentation, providing a set of object proposals. In the second stage, multi-scale crops are extracted from rendered views around each proposed instance and passed to CLIP to obtain object-level features. For our implementation, similar to above, we wish to establish an upper-bound of performance OpenMask3D can obtain. To that end, we skip Mask3D in the first stage and provide ground-truth 3D segmentation masks. We represent each instance with a

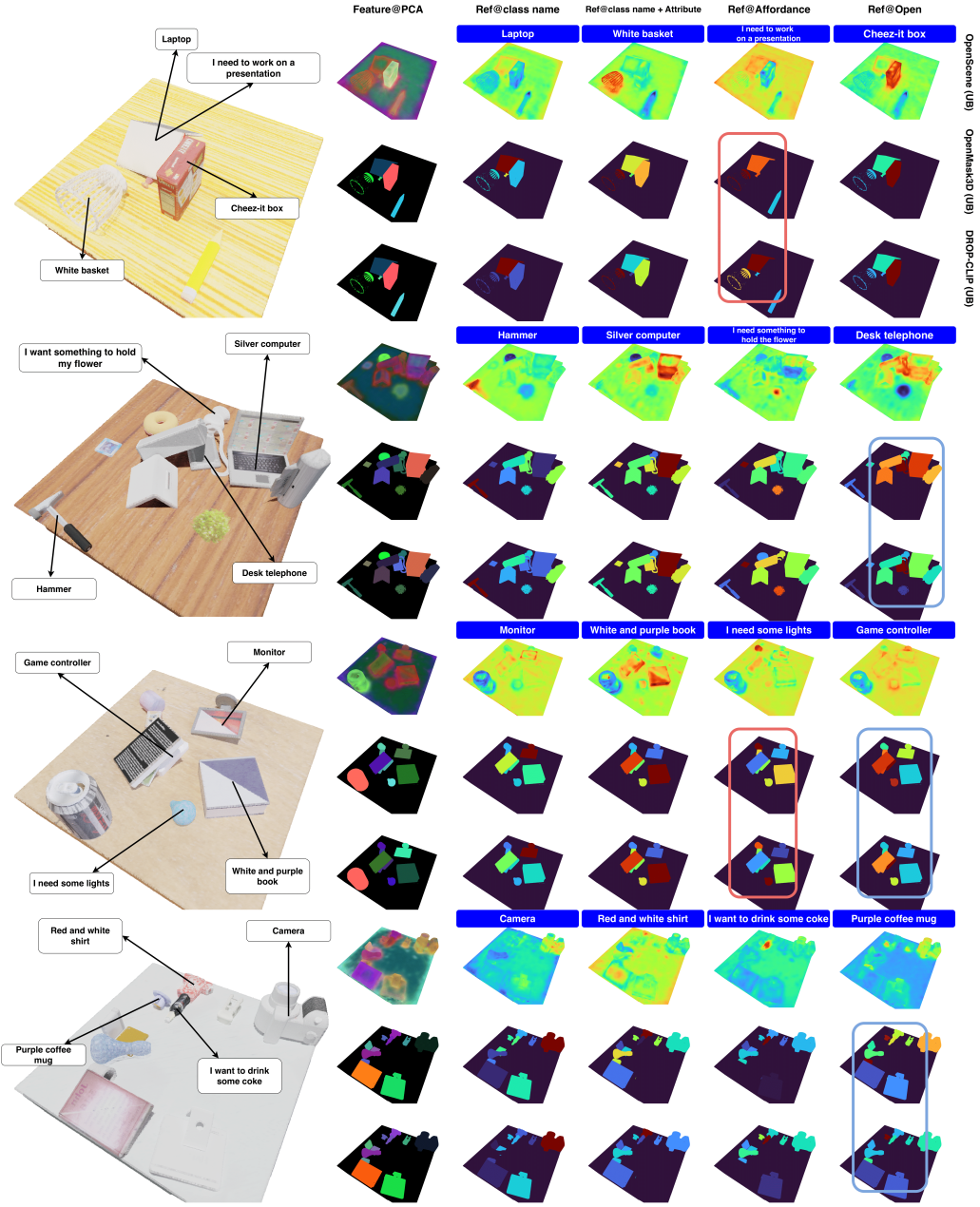


Figure 12: PCA feature and referring grounding visualization of baseline methods and DROP-CLIP. For each scene, we present results for OpenScene, OpenMask3D, and our DROP-CLIP (from top to bottom). The blue rectangle denotes cases where OpenMask3D suffers from distractor objects, while DROP-CLIP doesn't. The red rectangle denotes cases where OpenMask3D totally fails to ground the target, while DROP-CLIP succeeds.

pooled CLIP feature from 3 multi-scale crops of 0.1 expansion ratio, obtained through all of our 73 views and weighted according to the visibility map $\Lambda_{v,n}$ (see Sec. 3.2), as in the original paper.

A.5 Qualitative Results

We present qualitative results in several aspects to illustrate (1) How the object-centric priors help in multi-view feature fusion (Section A.5.1); (2) How do the distilled 3D features perform from single-view setting in MV-TOD semantic/referring segmentation tasks? (Section A.5.2).

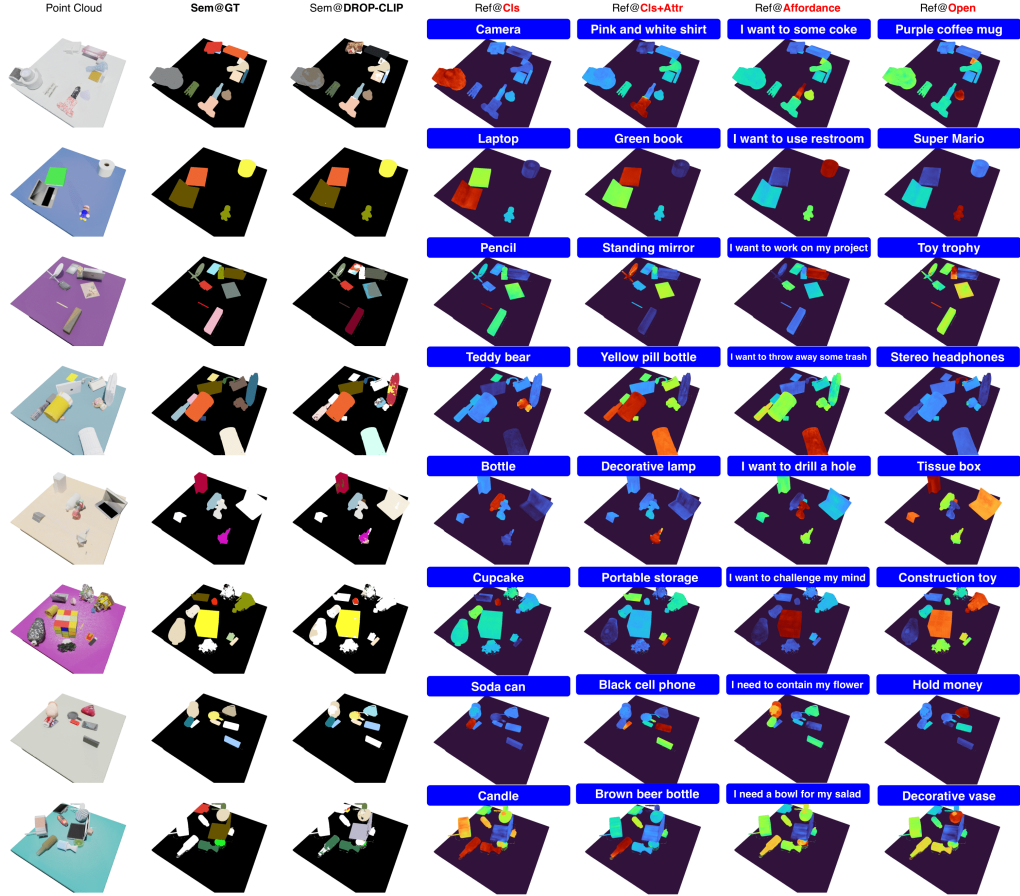


Figure 14: Semantic/Referring segmentation with our DROP-CLIP. In the **Sem@** columns, the same colors denote the same object category. The white parts mean that this part of the object is not activated by the corresponding class name query.

semantic informativeness metric in feature fusion, we add extra annotation in Fig. 12 and Fig. 13. The blue rectangle denotes the cases where OpenMask3D suffers from the distractors (i.e. multiple objects have high similarity score with the given query), while our DROP-CLIP is not. The red rectangle denotes the cases where OpenMask3D totally failed to ground the correct object, while our DROP-CLIP succeed. In conclusion, introducing semantic informativeness results in more robust object-level embeddings that in turn lead to higher grounding accuracy.

A.5.2 Referring / Semantic Segmentation Qualitative Results

Referring segmentation Since DROP-CLIP is not trained on closed-set vocabulary dataset but rather to reconstruct the fused multi-view CLIP features, the distilled features naturally live in CLIP text space. As a result, we can conduct referring expression segmentation in 3D with open vocabularies. We demonstrate this ability in Fig. 14 by showing the grounding results of the trained DROP-CLIP queried with different language expression, including *class name*, *class name + attribute*, *affordance*, and *open* instance-specific queries.

Semantic segmentation We present semantic segmentation results of our DROP-CLIP in Fig. 14. The white parts in Fig. 14 mean that this part of the object is not activated by the corresponding class name query.

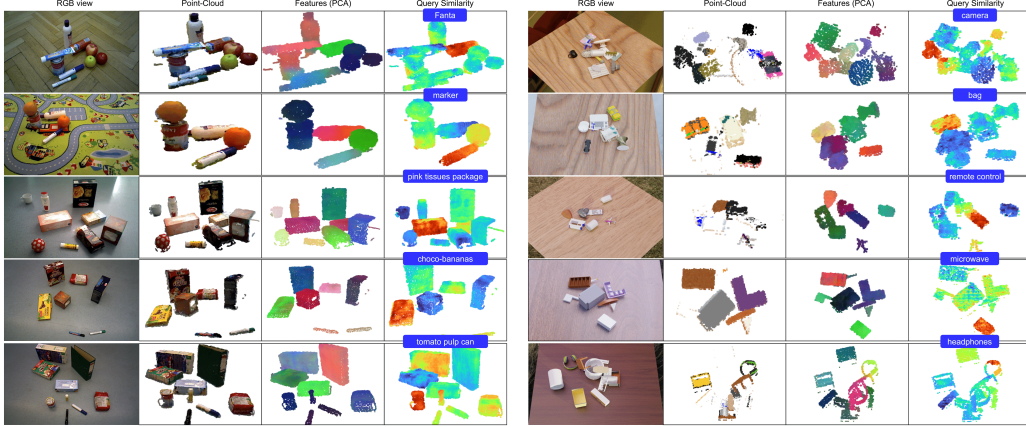


Figure 15: Visualization of referring segmentation examples in OCID-VLG ((left) and REGRAD (right) datasets.

A.6 Zero-Shot Transfer Experiments Details

To study the transferability of our learned 3D features in novel tabletop domains, in Sec. 4.3 we conducted single-view semantic segmentation experiments in the OCID-VLG dataset. In this setup, similar to our single-view MV-TOD experiments, we project the input RGB-D image to obtain a partial point-cloud and feed it to DROP-CLIP to reconstruct 3D CLIP features. To represent the point-clouds in the same scale as our MV-TOD training scenes, we sweep over multiple scaling factors and report the ones with the best recorded performance. For 2D baselines, the $mIoU$ and $mAcc@X$ metrics were computed based on the ground-truth 2D instance segmentation masks of each scene, after projected to 3D with the depth image and camera intrinsics and transformed to world frame, fixed at the center of the tabletop of each dataset.

A.6.1 Zero-Shot Referring Segmentation Experiments

Since methods LSeg and OpenSeg were fine-tuned for semantic segmentation, they are not suitable for grounding arbitrary referring expressions, but only category names as queries, which is why we conducted semantic segmentation experiments in our main paper. To further study zero-shot referring segmentation generalization, we conducted additional experiments in both OCID-VLG [31] and REGRAD [30] datasets. We compare with the MaskCLIP baseline, projected to 3D similar to above. For both the MaskCLIP baseline and DROP-CLIP, we use thresholding inference strategy, sweep over thresholds $\{0.4, \dots, 0.9\}$ and record the best configuration for both methods. We analyze the utilised datasets below:

OCID-VLG [31] connects 4-DoF grasp annotations from OCID-Grasp [74] dataset with single-view RGB scene images and language data generated automatically with templated referring expressions. We evaluate in one referring expression per scene for a total of 490 scenes, 165 from the validation and 325 from the test set of the *unique* split provided by the authors. We use the dataset’s referring expressions from *name* type as queries, after parsing out the verb (i.e. “*pick the*”), which contains open descriptions for 58 unique object instances, incl. concepts such as brand, flavor etc. (e.g. “*Kleenex tissues*”, “*Choco Krispies corn flakes*”, “*Colgate*”). For removing the table points, we use the provided ground-truth 2D segmentation mask to project only instance points to 3D for DROP-CLIP. We sweep over scaling factors $\{8, \dots, 16\}$ in the validation set and report the best obtained results.

REGRAD [30] focuses on 6-DoF grasp annotations and manipulation relations for cluttered tabletop scenes. Scenes are rendered from a pool of 50k unique ShapeNet [37] 3D models from 55 categories with 9 RGB-D views from a fixed height. We test in 1000 random scenes from *seen-val* split, using ShapeNet category names as queries. We note that as REGRAD doesn’t focus on semantics but

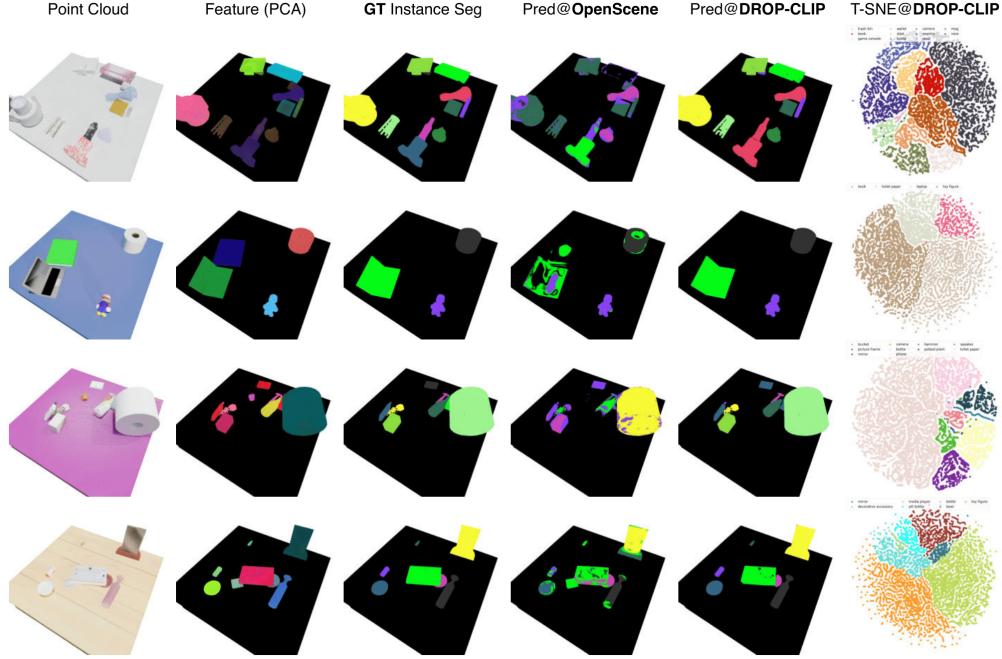


Figure 16: Zero-shot instance segmentation with our DROP-CLIP. In the *GT* and *Pred* columns, the same colors denote the same instance.

grasping, most of its objects are not typical household objects, but furniture objects (e.g. tables, benches, closets etc.) scaled down and placed in the tabletop. We filter out queries with such object instances and experiment with the remaining 16 categories that represent household objects (e.g. “bottle”, “mug”, “camera” etc.). We sweep over scaling factors $\{6, \dots, 20\}$. We use the filtered full point-clouds provided by the authors to identify the table points and remove them from each view.

More qualitative results for both datasets are illustrated in Fig. 15, while comparative results with MaskCLIP are given in Table 11. We observe that our method provides a significant performance boost across both domains (5.8% *mIoU* delta in OCID-VLG and 25.9% in REGRAD), especially in REGRAD scenes, where objects mostly miss fine texture and have plain colors, thus leading to poor MaskCLIP predictions compared to DROP-CLIP, which considers the 3D geometry of the scene. Failures were observed in cases of very unique referring queries in OCID-VLG (e.g. “Keh package”) and in cases of very heavily occluded object instances in both datasets.

Method	OCID-VLG		REGRAD	
	mIoU	Pr@25	mIoU	Pr@25
MaskCLIP $\rightarrow 3D$	40.4	45.2	33.2	39.0
DROP-CLIP	46.2	48.9	59.1	63.0

Table 11: Zero-shot referring segmentation results in OCID-VLG and REGRAD datasets

A.6.2 Zero-Shot 3D Instance Segmentation Experiments

Integrating spatial object priors via segmentation masks when fusing multi-view features grants separability in the embedding space. To illustrate that, we conducted zero-shot instance segmentation with our DROP-CLIP by directly applying DBSCAN clustering in the output feature space. In our experiments, we use the vanilla implementation of DBSCAN from `scikit-learn` package and set $\epsilon = 0.01$, `min_samples` = 2 for DROP-CLIP and $\epsilon = 0.01$, `min_samples` = 276 for OpenScene respectively. We observed that the points that belong to the same object instance are close to each other in the feature space, while significantly differ from the points that belong to other instances. We visualize several examples in Fig. 16, where we also conduct a t-SNE visualization to demonstrate the instance-level separability in the DROP-CLIP feature space.

A.6.3 Comparisons with SfM methods

In this section we compare DROP with modern 2D→3D feature distillation methods based on *Structure-from-Motion* (SfM), obtained via training NeRFs [13, 19, 16, 18] or 3D Gaussian Splatting (3DGS) [24, 75, 76, 77]. We highlight however that this is not really an “apples to apples” comparison, since SfM approaches differ from our method in philosophy and scope of application. In particular, SfM approaches perform **online** distillation in specific scenes, and thus require multiple camera images to distill, as well as significant time to do training / inference.

The obtained scene representation cannot be applied in new scenes, for which a new multi-view images dataset has to be constructed and a new NeRF / 3DGS be trained from scratch. In contrast, our approach relies on depth sensors to acquire 3D and does not need SfM reconstruction. The feature distillation is performed **offline once**, in the MV-TOD dataset, and thus can be applied zero-shot in novel scenes. Further, it does not require multiple camera images (works from single-view), does not require training and supports real-time inference. Nevertheless, we want to quantify the relative performance of DROP-CLIP with SfM methods that have been distilled for specific scenes.

We replicate the setup of the localization task from LERF [13] and the semantic segmentation task from LangSplat [24] for the ‘teatime’ scene of the LERF dataset. Results are presented in Table 12, where numbers for representative baselines LSeg [10], LERF [13], LangSplat [24] and Semantic Gaussians [75] are taken from corresponding papers. To signify the aforementioned differences in scope, in our table we also report number of views and training time required to obtain the representation (converted in v100 hours from time reported in corresponding papers) and whether / which external segmentation models (e.g. SAM [43]) is needed during test-time to deal with the ‘patchyness’ issue. The above demonstrate the practical benefits of our approach compared to SfM methods, as mentioned before, working from single-view, real-time performance, zero-shot application and no need for external segmentors. Regarding test results, we find that DROP scores lower to SfM baselines in both task variants, but significantly outperforms LSeg, which is the only other zero-shot baseline. The performance margin between DROP and object-centric 3DGS methods LangSplat and SemanticGaussians is significant, albeit the fact that these methods require SAM at test-time to inject the segmentation priors, whereas DROP doesn’t. This gap is justified when considering that our approach is zero-shot and didn’t have access to the 171 training scenes like the SfM baselines, as well as that the dataset queries are often referring to object parts (e.g. *hooves*, *bear nose* etc.), which DROP has not been designed for. Qualitative visualizations of DROP in the LERF scene are given in Fig. 17.

A.7 Robot Experiments

Setup Our robot setup consists of two UR5e arms with Robotiq 2F-140 grippers and an ASUS Xtion depth camera mounted from an elevated view between the arms. We conducted 50 trials in the Gazebo simulator [45] and 10 with a real robot. For simulation, we used 29 unique object instances from 9 categories (i.e. soda cans, fruit, bowls, juice boxes, milk boxes, bottles, cans, books and edible products). For real robot experiments, we mostly used packaged products and edibles. In each trial, we place 5-12 objects in a designated workspace area. Objects are either scattered across the workspace, packed together in the center or partially placed in the same area in order to emulate different levels of clutter. We provide a query indicating a target object using either category name, color/material/state attribute, user affordance (e.g. “I’m thirsty”), or open instance-level description, typically referring to the object’s brand (e.g. “Pepsi”, “Fanta” etc.) or flavor (e.g. “strawberry juice”, “mango juice” etc.) We note that distractor object instances of the same category as the target object are included in trials where query is not the category name.

Method	Modality	Num. Views	Train Time	Segm. Model	Results	
					Loc.	Sem.Segm.
LERF	SfM	171	112.5 min.	-	84.8	45.0
LangSplat	SfM	171	37.5 min.	SAM	88.1	65.1
SemanticGaussians	SfM	171	>2 hrs.	SAM	89.8	-
LSeg	RGB	1	0	-	33.9	21.7
DROP-CLIP	RGB-D	1	0	-	66.1	39.1

Table 12: Localization accuracy (%) and 3D semantic segmentation mIoU (%) on the ‘teatime’ scene of LERF dataset. We report number of views, training time and whether / which external models are needed to obtain the representation. Training times are converted to v100 hours from reported numbers in corresponding papers.

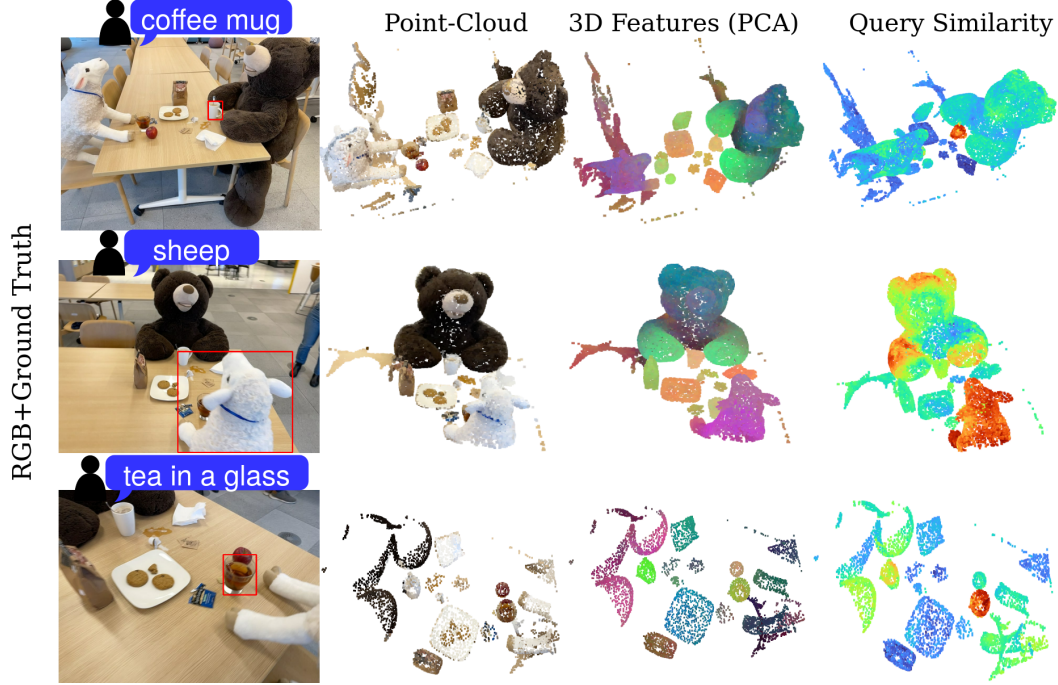


Figure 17: Visualizations of partial point-clouds, 3D DROP-CLIP features (PCA) and similarity heatmaps for three different queries in the ‘teatime’ scene of LERF dataset.

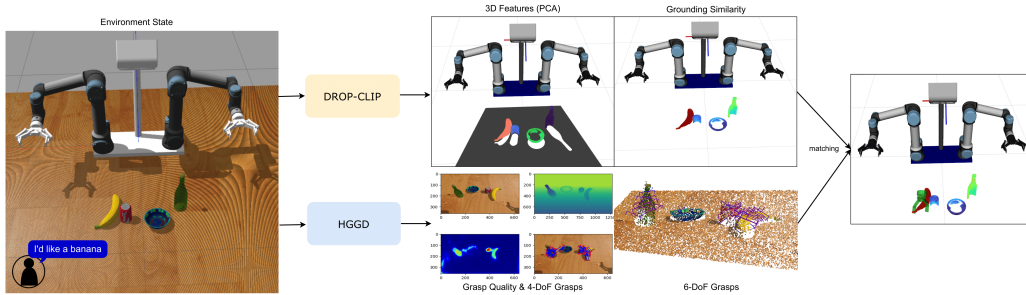


Figure 18: Illustration of robot system for language-guided 6-DoF grasping, using our DROP-CLIP for grounding (*top*), and HGGD network [44] for grasp detection (*bottom*).

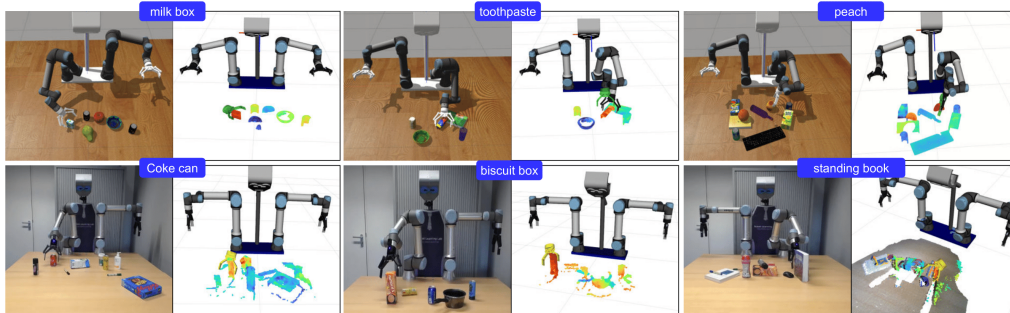


Figure 19: Visualization of robot experiments in Gazebo (*top*) and with a real robot (*bottom*).

Implementation We develop our language-guided grasping behavior in ROS, using DROP-CLIP for grounding the user’s query and RGB-D grasp detection network, HGGD [44], for generating

6-DoF grasp proposals. Our pipeline is shown in Fig. 18. We process the raw sensor point-cloud with RANSAC from `open3d` library with distance threshold 0.1, `ransac_n=3` and 1000 iterations to segment out the table points, and then upscale to $\times 10$. We use in-scene category names as negative prompts and do inference with a threshold of 0.7. To match the grounded object points with grasp proposals, we transform predicted grasps to world frame and move their center at the gripper’s tip. We then calculate euclidean distances between the gripper’s tip and the thresholded prediction’s center. In real robot experiments, we run statistical outlier removal from `open3d` with `neighbor_size=25` and `std_ratio=2.0` to remove noisy points from the prediction’s center. The top-3 closest grasps are given as goal for an inverse kinematics motion planner. We manually mark grasp attempts as success/failure in real robot and leverage the simulator state to do it automatically in Gazebo. Visualizations of simulated / real robot trials are illustrated in Fig. 19, experiments with grounding different objects with fine-grained attributes in Fig. 20, while related videos are included as supplementary material.

A.8 Detailed Related Work

In this section we provide a more comprehensive overview of comparisons with related work.

Semantic priors for CLIP in 3D A line of works aim to learn 3D representations that are co-embedded in text space by leveraging textual data, typically with contrastive losses [78, 79, 80]. CG3D [81] aims to learn a multi-modal embedding space by applying contrastive loss on 3D features from point-clouds and corresponding multi-view image and textual data, while using prompt tuning to mitigate the 3D-image domain gap. Most above methods lead to a degradation in CLIP’s open-vocabulary capabilities due to the fine-tuning stages. In contrast, our work leverages textual data not for training but for guiding multi-view visual feature fusion, hence leaving the learned embedding space intact from CLIP pretraining.

Spatial priors for CLIP in 3D Several works propose to leverage spatial object-level information to guide CLIP feature computation in 3D scene understanding context. OpenMask3D [6] leverages a pretrained instance segmentation method to provide object proposals, and then extracts an object-level feature by fusing CLIP features from multi-scale crops. Similarly, OpenIns3D [82] generates object proposals and employs a Mask-Snap-Lookup module to utilize synthetic-scene images across multiple scales. In similar vein, works such as Open3DIS [14], OVIR-3D [83], SAM3D [84], MaskClustering [85] and SAI3D [86] leverage pretrained 2D models to generate 2D instance-wise masks, which are then back-projected onto the associated 3D point cloud. All above approaches are two-stage approaches that rely on the instance segmentation performance of the pretrained model in the first stage, thus suffering from cascading effects when segmentations are not accurate or well aligned across views. In contrast, our method leverages spatial priors *during* the multi-view feature



Figure 20: Visualization of grounding queries in real robot trials, for baseline method $\text{MaskCLIP} \rightarrow 3D$ (bottom) and our DROP-CLIP (top), where we ensemble the predictions of our method with the 2D baseline. DROP-CLIP produces more robust features (middle column) which lead to crispier segmentation (right column.)

fusion process, and then distills the final features with a point-cloud encoder, and therefore is a single-stage method that does not require object proposals at test time.

Offline 3D CLIP Feature Distillation OpenScene [12] distills OpenSeg [9] multi-view features with a point-cloud encoder, while follow-up work Open3DSG [15] extends to scene graph generation by further distilling object-pair representations from other vision-language foundation models [87] as graph edges. CLIP-FO3D [21] replaces OpenSeg pixel-wise features with multi-scale crops from CLIP to further enhance generalization. All above works use dense 2D features and fuse point-wise, thus suffering from ‘patchyness’ issue. Further, these works distill features using 3D room scan data [22, 1], which lack diverse object catalogs and do not have to deal with the effects of clutter in the multi-view fusion process, as we do with the introduction of MV-TOD.

Online 3D CLIP Feature Distillation LERF [13] replaces point-cloud encoders with neural fields, and distills multi-scale crop CLIP features into a continuous feature field that can provide features in any region of the input space. The authors deal with the ‘patchyness’ issue using DINO regularization. Similar works OpenNeRF [19] and F3RM [16] use MaskCLIP to extract features and avoid DINO-regularization. All above works make the assumption that all views are equally informative and rely on dense number of views at test-time to resolve the noise in the distilled features. A more recent line of works replace NeRFs with 3D Gaussian Splatting (3DGS) [88] to improve inference time and memory consumption and perform similar feature distillation from LSeg, OpenSeg or CLIP multi-scale crops. Similar to our work, some 3DGS approaches [24, 75, 76, 77] also exploit spatial priors (i.e. segmentation masks) to distill object-level CLIP features, but do not perform view selection based on semantics. Further, 3DGS approaches lie in the same general family of works as fields, i.e., online distillation in specific scenes, requiring multiple camera images and computational resources to work at test-time. In contrast, our method is distilled offline in MV-TOD to reconstruct semantically-informed, view-independent 3D features from single-view RGB-D inputs, can be applied zero-shot in novel scenes without training, and enables real-time inference.