

OVI-MAP: Open-Vocabulary Instance-Semantic Mapping

Anonymous CVPR submission

Paper ID 1144

Abstract

001 Incremental open-vocabulary 3D instance-semantic mapping is essential for autonomous agents operating in complex everyday environments. However, it remains challenging due to the need for robust instance segmentation, real-time processing, and flexible open-set reasoning. Existing methods often rely on the closed-set assumption or dense per-pixel language fusion, which limits scalability and temporal consistency. We introduce OVI-MAP that decouples instance reconstruction from semantic inference. We propose to build a class-agnostic 3D instance map that is incrementally constructed from RGB-D input, while semantic features are extracted only from a small set of automatically selected views using Vision-Language Models. This design enables stable instance tracking and zero-shot semantic labeling throughout online exploration. Our system operates in real time and outperforms state-of-the-art open-vocabulary mapping baselines on standard benchmarks. The source code will be made publicly available.

019 1. Introduction

020 3D semantic and instance mapping are foundational capabilities for embodied perception, supporting downstream 021 tasks like language-conditioned navigation, manipulation 022 and facilitating queryable scene understanding for both 023 AR/VR and robotics [1, 15, 17, 21, 44]. In indoor environments, 024 voxel-based scene representations – most commonly Truncated Signed Distance Fields (TSDFs) – have become 025 a standard choice due to their real-time fusion, robustness 026 to drift, and ability to maintain dense and consistent geometry 027 for planning and interaction [11, 31, 32]. Recent systems 028 further couple volumetric reconstruction with semantics, 029 yielding panoptic maps that are temporally consistent 030 and easy to query [28, 45, 49, 54, 55].

031 However, existing pipelines are predominantly *closed-set*: they assume a fixed semantic ontology, learn class- 032 dependent predictors, and store integer class labels per 033 representation unit (*e.g.*, voxel, point). Extending these 034 designs to *open-set* recognition is non-trivial for several 035 036 037

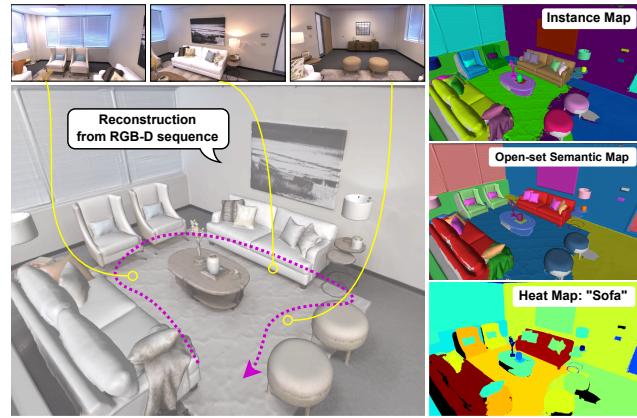


Figure 1. Overview of OVI-MAP. Given a streaming RGB-D sequence with camera poses, OVI-MAP incrementally reconstructs a volumetric 3D scene while maintaining a *class-agnostic* instance map. Semantic information is then assigned in a *zero-shot* manner using selectively chosen views, enabling open-set object recognition. Our method supports real-time, open-world scene reconstruction with instance-level semantic understanding.

038 reasons. First, open-set features extracted from Vision- 039 Language Models (VLMs) are high-dimensional and continuous; naively storing them at voxel-level resolution leads 040 to substantial computational and memory overhead.

041 Second, recent volumetric mapping systems [28, 49] rely 042 on semantic labels to guide instance segmentation and association. Without semantics, object instance grouping becomes 043 unstable and prone to fragmentation. Third, without 044 consistent 3D instances, aggregating per-pixel open-set 045 features over time is noisy due to occlusions, view-point 046 changes, background noise and inconsistent 2D segmentations. While recent segmentation models, such as 047 SAM [20], provide high-quality object proposals, running 048 them per frame is computationally expensive and therefore 049 unsuitable for real-time online mapping.

050 To address these limitations, we propose a pipeline 051 that decouples instance mapping from open-set semantic 052 recognition. Our approach first performs multi-view *class- 053 agnostic* instance segmentation and incrementally lifts these 054 instances into a global TSDF-based representation, produc- 055 ing a consistent *instance map* across frames without requir- 056 057 058

Method	3D Recon.	Inst. Map	Zero-Shot Sem.	Online
Ours	TSDF	✓	✓	✓
OVO-SLAM [27]	PCL/3DGS	✓	✓	✓
OpenFusion [49]	TSDF	✗	✗	✓
ConceptFusion [17]	PCL	✗	✓	✓
DCSEG [46]	3DGS	✓	✗	✗
OpenNeRF [8]	NeRF [29]	✗	✓	✗
Semantic Gauss [12]	3DGS	✗	✓	✗
OpenMask3D [43]		✗	✓	✗
OpenScene [33]		✗	✓	✗
Mask3D [40]		✗	✓	✗

Table 1. State-of-the-art online and offline semantic/instance mapping methods. PCL stands for Point Cloud, and 3DGS denotes 3D Gaussian Splatting [18].

ing semantic labels. On top of this instance map, we introduce an *object-centric view selection strategy* that identifies informative viewpoints for semantic extraction. Semantic features are queried from VLMs only when a new viewpoint provides non-redundant coverage of the 3D object’s surface, improving both efficiency and robustness while avoiding per-voxel storage of high-dimensional embeddings.

In summary, our contributions are:

- Class-agnostic online instance reconstruction pipeline that maintains a consistent 3D instance map independent of semantic categories.
- Object-centric incremental view selection mechanism that adaptively selects frames for feature extraction, reducing redundant VLM queries.
- An efficient, open-set semantic aggregation procedure that associates image regions with 3D instances and enables robust zero-shot semantic recognition.

We show on the ScanNet [6] and Replica [42] datasets that our proposed decoupled design leads to state-of-the-art instance mapping and open-set semantic segmentation while maintaining online processing rates.

2. Related Work

Open-Vocabulary Image Understanding has advanced rapidly with Vision-Language Models (VLMs), like CLIP [38] and SigLIP [56], that align visual and textual embeddings through large-scale contrastive pretraining. To enable pixel-level reasoning, many works produce dense or region-level semantic features. LSeg [22] predicts dense CLIP-aligned embeddings. OpenSeg [10] and OVSeg [26] refine segmentation using region proposals. SAM [20] and its derivatives [39, 47] offer category-agnostic masks, but produce over-segmented results. Works like Mask2Former [2, 23] provide semantic grounding only in closed-set, and X-Decoder [59, 60] provides a unified architecture for open-vocabulary segmentation. While these works show promising results on 2D images, they do not natively scale to 3D and there do not provide long-term, consistent segmentation across different views. We leverage these 2D models but focus on *incremental 3D* scenes and separate instance mapping from semantic recognition.

3D Semantic Understanding with VLMs. A growing class of methods lifts 2D open-vocabulary features into 3D by fusing multi-view semantic cues into volumetric, point-cloud, or neural field representations. OpenScene [33], PLA [7], and ConceptFusion [17] distill pixel-aligned VLM embeddings into point-level representations for open-vocabulary reasoning. Methods like CLIP-Fields [41], SemAbs [13], and 3DSS-VLG [48] propagate CLIP-aligned features across views and align them with text embeddings. More recently, Gaussian splatting and NeRF approaches [8, 12, 19, 25, 34, 37, 54, 58] process 3D scenes by directly embedding semantics into 3D Gaussian primitives or neural implicit fields. However, these methods often require global scene optimization or dense 3D semantic fields, making them sub-optimal for online, incremental mapping. In contrast, our method avoids this limitation and computes semantics only *per-instance* via selective view aggregation, yielding scalable real-time mapping.

3D Instance Segmentation and Panoptic Mapping is commonly performed directly on point clouds using architectures such as PointNet++ [35], MinkowskiNet [3], and Mask3D [40]. These models, however, rely on large-scale annotated 3D datasets and are trained in closed-set settings. To reduce reliance on manual labels, recent works [16, 51–53] propagate 2D mask priors into 3D, but still require frequent segmentation updates and do not target online operation. Incremental panoptic mapping in indoor scenes has been explored using TSDF-based fusion [11, 24, 28, 49, 57], yet these systems couple instance association with closed-set semantic label prediction [2, 14, 60], making them incompatible with open-set recognition.

Open-Set Instance-Level Semantic Mapping. OpenMask3D [43], Search3D [44] extracts CLIP-like features for 3D instances from multi-view crops, demonstrating that *instance-level* semantic aggregation is more stable than pixel-level fusion [33]. However, these works assume instance masks are available and do not address online mapping. Works like O3D-SIM [30] and DCSEG [46] combine clustering-based instance segmentation and class-dependent semantic masks from 2D VLMs, achieving offline instance-semantic segmentation with pre-defined semantic labels. Recent, related work OVO-SLAM [27] and OpenGS-SLAM [50] combine SLAM with open-vocabulary instance queries, but they still rely on time-consuming and over-segmented SAM [20] outputs and do not reason about view selection efficiency.

In summary, existing open-set 3D scene understanding methods either (1) store dense semantic fields in 3D, which is computationally prohibitive for incremental mapping, or (2) rely on offline or closed-set instance segmentation pipelines. Our work differs in two key aspects: (i) we reconstruct instances online in a class-agnostic manner, avoiding the dependence on closed-set semantic labels for

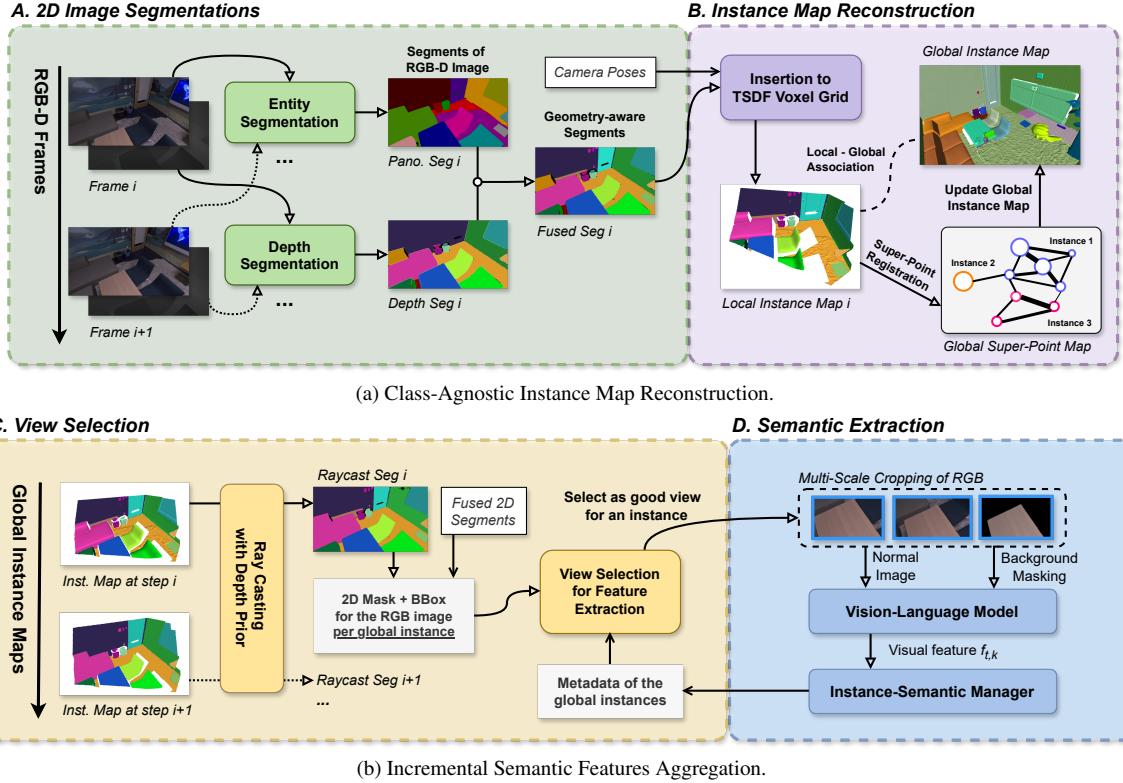


Figure 2. **Overview** of the proposed pipeline. **(a)** Class-agnostic instance map reconstruction: each RGB-D frame is segmented into entity proposals and refined with geometry-aware depth segmentation. The fused segments are lifted into 3D and incrementally integrated into a global TSDF-based instance map via super-point registration and spatial voting. **(b)** Incremental semantic feature aggregation: given the global instance map, each instance is re-projected into new frames via depth-guided ray casting. A view selection module identifies informative viewpoints based on object-centric coverage. Selected views are cropped at multiple scales and masked before being passed to a Vision-Language Model (VLM). The resulting features are aggregated per instance, producing stable open-set semantic embeddings.

object grouping, and (ii) we introduce *object-centric view selection* to efficiently aggregate VLM features, avoiding dense or redundant semantic fusion.

3. Proposed Method

Given a streaming RGB-D sequence $\{(I_t, D_t), \mathbf{T}_t\}_{t=1}^{\infty}$, where $\mathbf{T}_t \in \text{SE}(3)$ denotes the camera pose, I_t and D_t are the image and depth in the t -th frame, respectively, our goal is to incrementally construct: (i) a *class-agnostic 3D instance map*, and (ii) a *zero-shot semantic embedding* for each instance through selective multi-view feature extraction. The key idea is to decouple instance formation from semantic inference: instances are reconstructed purely from geometric and region-consistent evidence, and semantics are assigned only when *informative views* are observed. An overview of the pipeline is shown in Fig. 2.

3.1. Class-Agnostic Instance Map Reconstruction

We maintain the scene in a TSDF voxel grid \mathcal{V} [5], following [11], where each voxel v stores the values for TSDF:

$$(v_{\text{tsdf}}, v_{\text{weight}}, v_{\ell}),$$

with $v_{\ell} \in \mathbb{N}$ denoting the instance identity, or 0 if unassigned. Instance identities are represented via a dynamic set of 3D *super-points* $\mathcal{S} = \{S_1, \dots, S_K\}$, where each S_k corresponds to one globally consistent 3D instance and K denotes the current number of instances.

2D Entity Segmentation and Geometric Refinement. For each incoming RGB-D frame, we first compute a set of class-agnostic entity proposals as follows:

$$\{M_{t,j}\}_{j=1}^{N_t} = \text{EntitySeg}(I_t),$$

where $M_{t,j}$ is a binary mask identifying a coherent image region likely corresponding to one physical object. Following [36], the segmentation emphasizes object-level boundaries while ignoring category semantics.

To improve segmentation around contacts and clutter, we compute a geometric segmentation G_t , in the t -th RGB-D frame, based on depth discontinuities [9] and use it to improve the obtained 2D masks as follows:

$$\hat{M}_{t,j} = \text{MaskFusion}(M_{t,j}, G_t),$$

where MaskFusion further split mask G_t into smaller pieces according to $\{M_{t,j}\}_{j=1}^{N_t}$ to avoid under-segmentation. This

191 reduces over-merge errors where adjacent objects share texture
 192 or color, by using geometric discontinuities indicating
 193 boundaries not visible in appearance space.

194 **3D Lifting.** Each refined segment $\hat{M}_{t,j}$ is lifted into a 3D
 195 point cloud $P_{t,j}$ in the global coordinate frame using camera
 196 pose \mathbf{T}_t and depth projection as follows:

$$197 P_{t,j} = \{\mathbf{x} = \mathbf{T}_t \cdot (D_t(\mathbf{u})K^{-1} \begin{bmatrix} \mathbf{u} \\ 1 \end{bmatrix}) : \mathbf{u} \in \hat{M}_{t,j}\},$$

198 where K denotes the camera intrinsics, \mathbf{u} is a 2D point from
 199 the region defined by mask $\hat{M}_{t,j}$, $D_t(\mathbf{u})$ is the depth ob-
 200 served at \mathbf{u} , and \mathbf{x} is the 3D point. This produces a partial
 201 surface patch corresponding to a single candidate object.

202 **Instance Association via Spatial Voting.** To determine if
 203 $P_{t,j}$ corresponds to an existing instance, we examine which
 204 instance label (*i.e.*, index of a 3D instance) appear the most
 205 frequently in voxels covering the points in $P_{t,j}$.

206 Let $V(\mathbf{x})$ denote the voxel contains \mathbf{x} , we define the voting
 207 of the super-point S_k with label k as:

$$208 \Omega_{j,k} = |\{\mathbf{x} \in P_{t,j} : V(\mathbf{x})_k = k\}|,$$

209 where $V(\mathbf{x})_k$ is the instance label of voxel $V(\mathbf{x})$. The as-
 210 signment decision is the following:

$$211 k^* = \arg \max_k \Omega_{j,k}, \quad \text{if } \Omega_{j,k^*} > \theta_{\text{assoc}}, \text{ assign } P_{t,j} \mapsto S_{k^*};$$

212 where θ_{assoc} is a fixed overlapping threshold, otherwise a
 213 super-point S_{new} is created with a new label.

214 This association depends entirely on *spatial consistency*,
 215 not semantic similarity. Thus, unlike TSDF-based panoptic
 216 fusion pipelines [28, 57], we do not require a closed label
 217 set or hierarchical class priors.

218 **Incremental TSDF Fusion and Label Stabilization.** Once
 219 an instance with label k^* is selected for $P_{t,j}$, its surface
 220 points are integrated into the TSDF grid using standard
 221 weighted fusion [31]. Instance identity is reinforced via per-
 222 voxel label-support accumulation as follows:

$$223 O_v(k^*) \leftarrow O_v(k^*) + 1 \quad \forall v = V(\mathbf{x}), \mathbf{x} \in P_{t,j}.$$

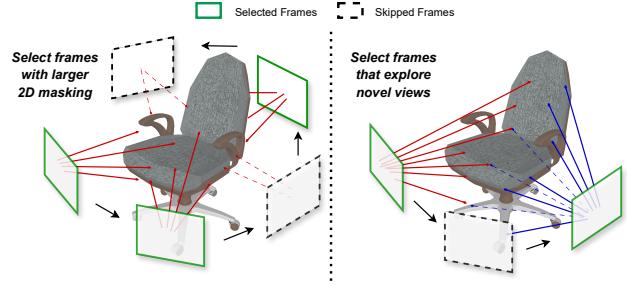
224 After processing frame t , each voxel updates its label as:

$$225 v_\ell = \arg \max_k O_v(k) \quad \forall v \in \mathcal{V}.$$

226 This forms a temporally stable, class-agnostic 3D instance
 227 map \mathcal{M} that progressively improves as the scene is observed
 228 from more viewpoints. Moreover, over-segmentation and
 229 multi-view fusion are handled by merging spatially close
 230 super-points, introduced in the supplementary materials.

231 3.2. View-Adaptive Semantic Feature Aggregation

232 Once object geometry is stable, we compute open-set de-
 233 scriptors. Rather than aggregating pixel-level VLM features



234 **Figure 3. Comparison of view selection strategies.** The left il-
 235 lustration shows the pixel-counting strategy [43], which prioritizes
 236 frames with larger object masking area, often leading to redundant
 237 front-facing views. The right illustration depicts our proposed
 238 *object-centric view coverage* method, which maintains a spherical
 239 coordinate map over bins, where if a bin contains 1, the
 240 instance has been observed from that direction.

241 across all frames – leading to noise, redundancy, and mem-
 242 ory inefficiency – we extract semantic features *only* when a
 243 view provides novel evidence about object appearance.

244 **Object-Centric View Coverage Representation.** For each
 245 instance S_k , we maintain its *view coverage* on the unit
 246 sphere. Let $\text{Cov}_k \in \{0, 1\}^{180 \times 240}$ denote the spherical
 247 coordinate map over bins, where if a bin contains 1, the
 248 instance has been observed from that direction.

249 We initialize S_k a centroid \mathbf{c}_k of its 3D bounding box
 250 when it gets an instance mask with area larger than a fixed
 251 threshold. For every visible point $\mathbf{x} \in P_{t,k}$ at frame t , we
 252 compute its viewing direction:

$$253 \mathbf{d}_{t,\mathbf{x}} = \frac{\mathbf{x} - \mathbf{c}_k}{\|\mathbf{x} - \mathbf{c}_k\|},$$

254 and convert $\mathbf{d}_{t,\mathbf{x}}$ into spherical coordinates (θ, ϕ) to com-
 255 pute the current view coverage in the map Cov_k .

256 The novelty of the current view is measured by the ratio:

$$257 \eta_{t,k} = \frac{|\text{BinsNewOccupied}(P_{t,k})|}{|\text{BinsOccupied}(P_{t,k})|}.$$

258 If $\eta_{t,k} > \theta_{\text{novel}}$, the view reveals new geometric or ap-
 259 pearance information and is selected for semantic extrac-
 260 tion, where θ_{novel} is a predefined threshold. After a view is
 261 selected, all bins in Cov_k occupied by (θ, ϕ) are set to 1.

262 As illustrated in Fig. 3, this object-centric criterion dif-
 263 fers fundamentally from the ‘visible pixel counting’ heuris-
 264 tic [43]: while the latter repeatedly selects large, front-
 265 facing views that fail to capture full object shape, our
 266 coverage-based selection explicitly favors novel viewpoints
 267 that expand the explored surface area, yielding more diverse
 268 and informative observations.

269 **Object-Level VLM Feature Extraction.** For the selected
 270 frame of an instance, we extract the visual embeddings from

two kinds of object crops: 1) a cropped image containing the object extent, and 2) a masked version where background pixels are removed:

$$\mathbf{f}_{t,k}^{(1)} = F_{\text{VLM}}(I_t, P_{t,k}), \quad \mathbf{f}_{t,k}^{(2)} = F_{\text{VLM}}(I_t \odot M_{t,k}).$$

We average these and update a running embedding:

$$\mathbf{f}_k \leftarrow \frac{w_{\text{sum}}}{w_k + w_{\text{sum}}} \mathbf{f}_k + \frac{w_k}{2(w_k + w_{\text{sum}})} \cdot (\mathbf{f}_{t,k}^{(1)} + \mathbf{f}_{t,k}^{(2)}),$$

where w_k denotes the number of visible object pixels in frame k and w_{sum} is the cumulative pixel count over all previous observations. This adaptive update produces a stable, view-invariant embedding that captures semantic consistency across frames, enabling robust zero-shot reasoning and language-conditioned retrieval.

Why Decoupling Works. Since instance mapping depends only on geometry and does not require consistent labeling, it remains stable and efficient in open-world environments. Semantic reasoning is then performed only when evidence is strong and informative, mitigating accumulated noise and reducing VLM calls significantly and controllably.

4. Experiments

Implementation Details. We adopt CropFormer [36] for 2D entity segmentation and the depth-based segmentor from [9] to obtain geometric boundaries. Our pipeline is implemented on top of Voxblox++ [11] using a TSDF [31] representation with voxel size of 0.1m, following [28] but without closed-set semantic dependencies. For semantic extraction, we use SigLIP [56] as the vision-language backbone to obtain zero-shot embeddings from RGB crops. For evaluation, text labels are encoded with the same model, and matched to instance features via cosine similarity. We run all methods on an Nvidia RTX3090 GPU and an Intel Core i7-12700K CPU.

Datasets. We evaluate our method on the *Replica* [42] and *ScanNet* [6] datasets, which provide dense RGB-D sequences with ground-truth camera poses, instance labels, and semantic annotations. For fair comparison, all methods are evaluated under identical input trajectories and reconstructed geometry. We run each sequence with 200 frames evenly sampled, as done in baselines [8, 12, 33]. For evaluation, we use the 51-class label set defined for Replica and the ScanNet200 label set for ScanNet.

4.1. Instance Segmentation Results

We first evaluate the quality of the reconstructed instance maps. Each reconstructed map is projected to the ground-truth mesh using a kNN-based nearest-neighbor mapping [4] with a fixed distance threshold, allowing per-vertex comparison with the annotated ground truth. We compare our approach with three state-of-the-art open-set methods: Mask3D [40], Segment3D [16], and OVO-SLAM [27].

	Method	Online	mIoU	AP ₇₅	AP ₅₀	AP ₂₅
Replica	Mask3D	✗	23.1	14.3	31.2	56.2
	Segment3D	✗	29.3	6.5	14.5	24.9
	OVO-SLAM	✓	42.7	11.1	23.6	32.8
	Ours	✓	<u>36.3</u>	22.0	50.8	76.7
ScanNet	Mask3D	✗	47.6	16.9	36.1	47.8
	Segment3D	✗	42.0	2.5	9.3	16.5
	OVO-SLAM	✓	39.8	2.0	7.4	14.4
	Ours	✓	41.2	9.8	<u>24.0</u>	<u>37.4</u>

Table 2. **Instance segmentation performance** on Replica [42] and ScanNet [6]. The *Online* column indicates whether the method performs incremental mapping during exploration. We achieve comparable or better accuracy to offline approaches while consistently outperforming prior online systems, particularly at high IoU thresholds (AP_{75}). Mask3D was trained on ScanNet.

Mask3D and Segment3D are offline volumetric networks predicting per-vertex instance masks, whereas OVO-SLAM and our method perform online, incremental mapping.

Metrics. We report mean IoU (mIoU) and instance-level average precision (AP) at IoU thresholds of 25%, 50%, and 75% following [16]. Best and second-best scores are highlighted. Offline methods are evaluated on the reconstructed meshes mapped to the ground truth ones.

Results. As shown in Table 2, our method achieves similar or higher instance-level precision compared to offline networks, while maintaining online operation. In particular, it surpasses the recent OVO-SLAM across all thresholds, with a substantial margin in AP_{75} , demonstrating the robustness of our class-agnostic instance tracking and fusion. It is important to note that, while Mask3D performs the best on ScanNet, it was trained on that dataset.

4.2. Semantic Segmentation Results

We next evaluate the open-vocabulary semantic performance. To further assess real-time capability, we benchmark both our OVI-MAP and the baseline OVO-SLAM [27] under a *30 FPS constraint*, where only every n -th frame is processed for semantic extraction to run in real time. On the Replica dataset, semantics are computed every 10th frame for our method and every 50th frame for OVO-SLAM, while on ScanNet, we process every 10th and 30th frame, respectively. Thanks to our lightweight and view-selective semantic extraction pipeline, OVI-MAP can perform semantic updates more frequently without exceeding real-time limits.

Baselines. We compare with offline and online open-vocabulary mapping systems. Among the offline methods, Mask3D+OpenMask3D (OM3D) and Segment3D+OM3D combine an instance segmentation backbone with zero-shot semantic labeling: Mask3D [40] and Segment3D [16] generate 3D instance masks from the mesh input, while OpenMask3D [43] assigns open-vocabulary semantic feature to

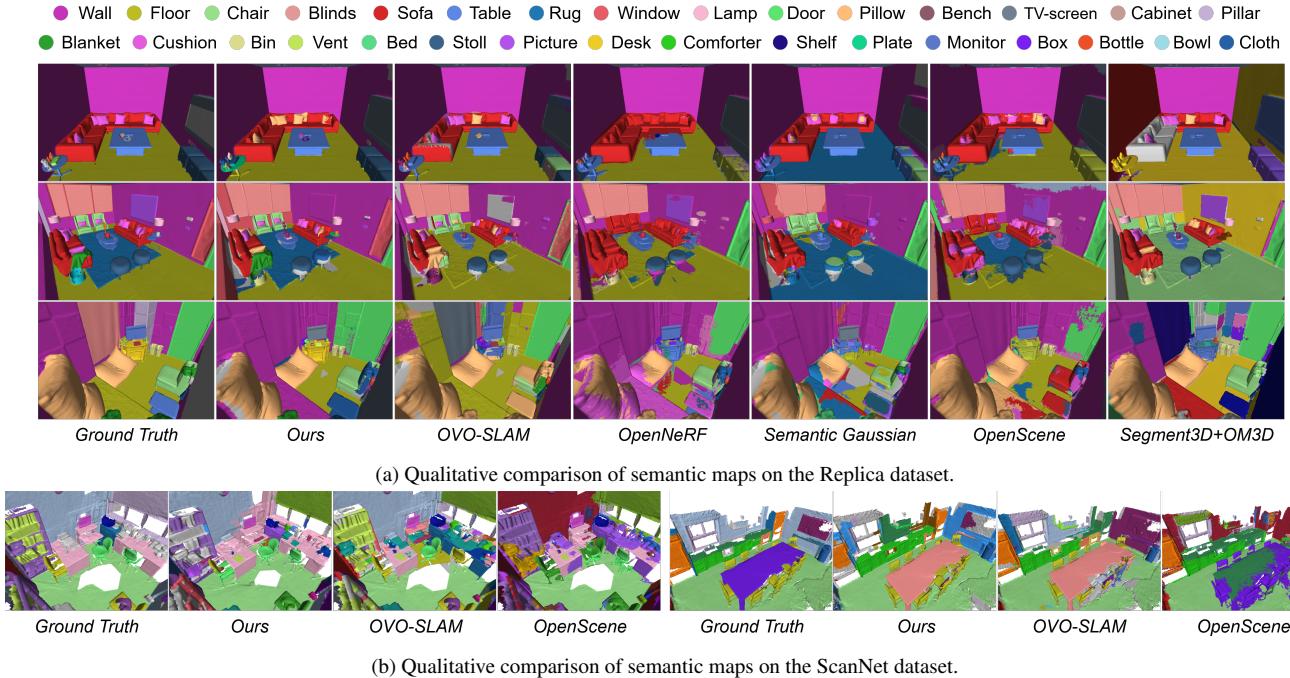


Figure 4. **Open-vocabulary 3D semantic maps** aligned to the ground-truth label sets of the respective datasets. We compare our method with online [27] and offline [8, 12, 16, 33, 43] approaches on the *Replica* (a) and *ScanNet* (b) datasets. Colors correspond to the semantic classes defined in each dataset. Gray regions indicate unobserved areas. OVI-MAP produces spatially coherent and semantically accurate reconstructions, maintaining sharp instance boundaries and consistent semantics throughout incremental mapping. Minor discrepancies in color (e.g., pillows in (a) and the table in (b)) arise from closed-set label mapping, where the ground truth uses alternative class names such as “cushion” or “dining table”.

	Method	Online	mIoU	mAcc	AP ₂₅	AP ₅₀	AP _{all}
Replica	Mask3D+OM3D	✗	18.7	29.4	6.8	4.5	3.2
	Segment3D+OM3D	✗	17.3	30.9	19.9	15.2	9.0
	Semantic Gaussian	✗	16.1	22.1	—	—	—
	OpenScene	✗	19.8	33.9	—	—	—
	OpenNeRF	✗	19.2	30.9	—	—	—
	OVO-SLAM	✓	24.9	34.0	28.1	17.5	9.1
	Ours	✓	26.5	32.2	34.5	21.2	8.5
	OVO-SLAM (30 fps)	✓	21.8	27.5	21.5	15.2	8.1
ScanNet	Ours (30 fps)	✓	27.0	32.5	31.8	17.7	8.0
	Mask3D+OM3D	✗	8.6	17.5	10.4	8.0	5.1
	Segment3D+OM3D	✗	4.5	13.5	5.0	3.9	2.3
	Semantic Gaussian	✗	10.7	17.8	—	—	—
	OpenScene	✗	10.3	17.3	—	—	—
	OpenNeRF	✗	Fail	Fail	—	—	—
	OVO-SLAM	✓	14.6	27.8	19.4	12.6	5.5
	Ours	✓	17.5	27.6	23.4	15.7	7.2
	OVO-SLAM (30 fps)	✓	15.5	26.9	19.4	13.7	5.8
	Ours (30 fps)	✓	<u>16.3</u>	25.4	<u>21.1</u>	<u>14.4</u>	7.0

Table 3. Comparison of **open-vocabulary 3D semantic and instance segmentation** performance on Replica [42] and ScanNet [6]. OM3D denotes OpenMask3D [43]. The 30 FPS rows correspond to experiments where only every n -th frame is processed to match real-time performance. Our method achieves the best open-vocabulary semantic-instance accuracy among online systems, and maintains performance comparable to offline approaches even under real-time constraints.

each reconstructed instance with the RGB-D input. We also include Semantic Gaussian [12], OpenScene [33], and OpenNeRF [8], which perform per-vertex open-vocabulary

reasoning directly from scene representations without explicit instance segmentation. For online mapping, we compare to OVO-SLAM [27].

Metrics. We report mean IoU (mIoU) and accuracy (mAcc) following [8], as well as instance-level precision at different IoU thresholds (AP_{25} , AP_{50} , AP_{all}) as done in [43]. Since the ground-truth annotations in both datasets are provided as closed-set labels, open-vocabulary features are projected onto the label set defined by the datasets, following [8]. The semantic feature of each reconstructed instance is compared against the SigLIP-encoded text embeddings of these class names, and the most similar label is assigned based on cosine similarity to enable consistent quantitative evaluation.

Results. As shown in Table 3, our method achieves the highest instance-level precision among the online systems and remains competitive with, or superior to, offline pipelines such as OpenScene and OpenNeRF. On both Replica and ScanNet, our approach improves the semantic-instance AP metrics by a large margin over OVO-SLAM, particularly at higher IoU thresholds (AP_{50} and AP_{25}), demonstrating that our view-selection-based semantic aggregation yields more reliable and consistent semantics. Under the 30 FPS real-time constraint, OVI-MAP retains most of its original performance, achieving only marginal



Figure 5. **Instance highlighting from arbitrary text queries.** Given natural language prompts, our system retrieves and highlights corresponding 3D instances based on the learned vision-language embeddings. The examples demonstrate zero-shot grounding of both concrete (“pillow”, “toilet”) and abstract (“where to sleep”, “where is the music”) concepts in reconstructed scenes. Darker tones indicate higher cosine similarity between an object and the query.

	Method	mIoU	mAcc	AP₂₅	AP₅₀	AQ_↓
Replica	Random 8 Views	23.8	30.4	31.4	18.6	50.3
	Pixel Counting	26.5	31.8	33.2	19.8	18.7
	View Coverage (Ours)	26.5	32.2	34.5	21.2	8.6
	<i>GT Inst. + Pixel Cnt.</i>	38.0	50.3	37.6	37.6	20.7
	<i>GT Inst. + View Cov.</i>	37.6	48.4	36.2	36.2	10.4
ScanNet	Random 8 Views	14.7	22.4	19.2	13.7	23.0
	Pixel Counting	17.5	27.5	24.7	17.0	15.8
	View Coverage (Ours)	17.5	27.2	23.4	15.7	7.6
	<i>GT Inst. + Pixel Cnt.</i>	34.8	48.9	39.8	39.8	23.0
	<i>GT Inst. + View Cov.</i>	30.7	44.8	34.7	34.7	12.3

Table 4. **View selection strategies** for semantic feature extraction: (1) random view selection, (2) larger visible-pixel counting similar to [43], and (3) our proposed *object-centric view coverage* method. *AQ* denotes the average number of VLM queries per instance. Our approach achieves comparable zero-shot instance-level semantic precision (AP_{25}/AP_{50}) to pixel-counting while requiring substantially fewer (47%) VLM queries, enabling efficient real-time operation. The *GT Inst. + Pixel Cnt.* and *GT Inst. + View Cov.* configurations provide an upper bound using ground-truth instance masks.

375 drops in *mIoU* and *AP* scores, whereas OVO-SLAM exhibits a substantial degradation on Replica. Our results out-
376 perform offline methods even in the real-time regime.
377

378 Figure 4 presents qualitative results of the reconstructed
379 semantic maps on Replica and ScanNet. Our method pro-
380 duces spatially coherent and semantically rich reconstruc-
381 tions, maintaining clear object boundaries and accurate la-
382 bel assignments across diverse viewpoints.

383 **Open Query Searching.** Figure 5 shows the map high-
384 lighting the most related objects in the scene with the given
385 text queries. These results also demonstrate the accurate
386 instance segmentation accuracy of our method.

387 4.3. Ablation Studies

388 We conduct ablation studies to analyze the effect of the pro-
389 posed components, with a focus on the view-selection strat-
390 egy for semantic feature aggregation and the trade-off be-

tween accuracy and computational efficiency.

View Selection Strategy. We evaluate the influence of different view selection strategies on semantic feature aggregation, focusing on the trade-off between recognition accuracy and computational efficiency. The goal is to minimize redundant queries to the vision-language model (VLM) while maintaining accurate open-vocabulary semantics.

We compare three strategies for selecting views per instance: (1) *Random* selection of eight frames per instance, serving as a naive baseline; (2) *Pixel-counting* selection, which prioritizes frames with the largest visible mask area following OpenMask3D [43]; and (3) our proposed *object-centric view coverage* selection, which measures how much of an object’s surface has been observed and selects frames that contribute novel surface regions.

As shown in Table 4, random view selection, as expected, yields poor performance. Our object-centric view coverage strategy achieves higher accuracy than pixel-counting, while using less than half the number of VLM queries (47.0% on average). This efficiency gain arises from explicitly modeling viewpoint novelty, which avoids redundant observations while preserving semantic consistency.

Incremental Performance. Fig. 6 shows how semantic accuracy evolves as more views are observed. Both our *view coverage* and *pixel-counting* strategies exhibit consistent improvement as additional frames are processed, demonstrating the benefits of multi-view semantic aggregation. Our view coverage approach achieves slightly better accuracy while requiring significantly fewer VLM queries per instance, as shown in the bottom-right plot. In contrast, post-hoc semantic extraction methods – such as OpenMask3D, Segment3D, and OpenScene [33] – only converge after the full sequence is processed, offering no incremental feedback during exploration. The results confirm that our coverage-based selection enables efficient, real-time semantic reasoning that scales with scene exploration, effectively

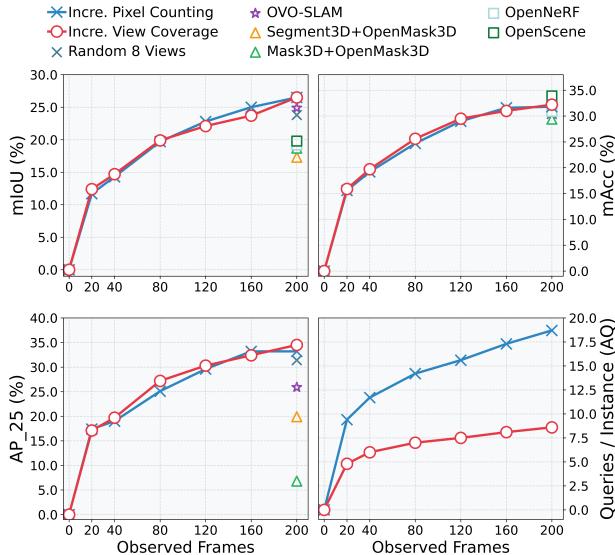


Figure 6. **Incremental performance** on Replica as the number of observed views increases. We compare our incremental semantic aggregation using the proposed *view coverage* strategy against pixel-counting and other baselines. Our method achieves comparable semantic accuracy (mIoU, mAcc, AP_{25}) while requiring significantly fewer VLM queries per instance, demonstrating efficient and scalable open-vocabulary mapping during online exploration.

427 bridging the gap between accuracy and online performance.

428 **Influence of 2D Instance Segmentation.** We further study
429 how the quality of instance segmentation affects down-
430 stream semantic feature extraction. Table 5 compares dif-
431 ferent methods for 2D instance segmentation: (1) *SAM2*,
432 which provides high-recall but noisy over-segmentations;
433 (2) *CropFormer* [36], our chosen segmentor, producing
434 compact and spatially coherent entities; and (3) *GT 2D in-*
435 *stance masks*, serving as an oracle reference.

436 The results show that improving instance segmentation
437 quality directly enhances semantic accuracy. CropFormer
438 outperforms SAM2 in both instance- and semantic-level
439 metrics, yielding a +5.7 improvement in AP_{50} and a +6.8
440 gain in semantic mIoU. This confirms that cleaner and more
441 consistent instance boundaries lead to more discriminative
442 VLM embeddings. The use of ground-truth instance masks
443 further boosts all metrics, demonstrating the upper bound
444 achievable with perfect segmentation.

445 **Feature Fusion.** We analyze how semantic features from
446 multiple selected views should be fused to form stable
447 instance-level embeddings. Table 6 compares several strate-
448 gies: simple *averaging*, *weighted averaging* by the number
449 of visible object pixels, and two *feature clustering* variants
450 based on cosine similarity and L_1 distance.

451 Averaging provides a strong baseline, confirming that
452 multi-view aggregation mitigates noise in individual pre-

Method	Instance Segmentation			Semantic Segmentation			
	mIoU	AP ₇₅	AP ₅₀	mIoU	mAcc	AP ₂₅	AP ₅₀
SAM2	36.6	14.2	27.8	20.2	26.4	26.2	18.6
Cropformer (Ours)	36.3	22.0	50.8	26.9	33.2	36.4	22.0
GT 2D Inst. Seg.	53.9	38.8	65.7	26.6	33.8	36.1	27.3
<i>GT Instances Map</i>	100	100	100	38.0	50.3	37.6	37.6

Table 5. **Influence of instance segmentation:** instance map built with 2D instance masks from (1) SAM2, (2) CropFormer, and (3) ground-truth 2D instance masks. All methods use the same view-selection strategy for semantic fusion. CropFormer yields the highest zero-shot semantic precision, while the *GT Instances* row represents an upper bound with perfect instance masks.

Method	mIoU	mAcc	AP ₂₅	AP ₅₀	AP _{all}
Averaging	26.5	31.8	<u>33.2</u>	<u>19.8</u>	8.3
+ Weighted by Pixel Count	26.9	33.2	36.4	22.0	8.6
Cluster - Max. Cosine Sim.	25.1	31.9	33.1	19.7	8.4
Cluster - Min. L_1 Distance	24.8	31.2	32.8	19.5	8.3

Table 6. **Semantic feature fusion by** (1) averaging, (2) weighted averaging by normalized pixel counts (Ours), (3) selecting the feature with maximum average cosine similarity to others, and (4) selecting the feature with minimum average L_1 distance to others.

453 diction. Weighting by visible pixel count yields the best
454 results across all metrics, as it naturally emphasizes clearer
455 and better-framed object observations. Clustering-based
456 methods offer no additional benefit and sometimes degrade
457 performance due to overemphasis on redundant features.

458 These results indicate that a simple yet visibility-aware
459 fusion is sufficient for robust zero-shot semantic reasoning,
460 avoiding the complexity of feature clustering.

461 5. Conclusions

462 Our work OVI-MAP demonstrates that open-vocabulary se-
463 mantic mapping can be performed incrementally and effi-
464 ciently without sacrificing accuracy. By combining class-
465 agnostic instance reconstruction with language-guided fea-
466 ture aggregation, we show that reliable semantic reasoning
467 is possible even under online constraints. The proposed
468 object-centric view coverage strategy is central to this ca-
469 pability, enabling informed view selection that maintains
470 semantic consistency while substantially reducing VLM
471 queries. Together, these components bridge the gap be-
472 tween offline open-set understanding and real-time 3D per-
473 ception. The source code will be publicly released.

474 **Limitations.** Our method still depends on the quality of 2D
475 segmentation [36], which limits performance on small or
476 visually complex objects. Semantic embeddings extracted
477 from masked RGB crops are also affected by segmen-
478 tation errors and background bias. Finally, current vision-
479 language models provide only weak alignment between vi-
480 sual and textual spaces, causing ambiguous label assign-
481 ments. Future research will explore tighter vision-language
482 coupling and more adaptive feature fusion to enhance ro-
483 bustness in cluttered, real-world environments.

484

References

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

- [1] Boyuan Chen, Fei Xia, Brian Ichter, Kanishka Rao, Keerthana Gopalakrishnan, Michael S Ryoo, Austin Stone, and Daniel Kappler. Open-vocabulary queryable scene representations for real world planning. In *IEEE Int. Conf. on Robotics and Automation*, pages 11509–11522, 2023. 1
- [2] Bowen Cheng, Ishan Misra, Alexander G Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1290–1299, 2022. 2
- [3] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3075–3084, 2019. 2
- [4] Thomas Cover and Peter Hart. Nearest neighbor pattern classification. *IEEE transactions on information theory*, 13(1): 21–27, 1967. 5
- [5] Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *Conference on Computer Graphics and Interactive Techniques*, pages 303–312, 1996. 3
- [6] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5828–5839, 2017. 2, 5, 6, 3
- [7] Runyu Ding, Jihan Yang, Chuhui Xue, Wenqing Zhang, Song Bai, and Xiaojuan Qi. Pla: Language-driven open-vocabulary 3d scene understanding. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 7010–7019, 2023. 2
- [8] Francis Engelmann, Fabian Manhardt, Michael Niemeyer, Keisuke Tateno, Marc Pollefeys, and Federico Tombari. Opennerf: open set 3d neural scene segmentation with pixel-wise features and rendered novel views. In *Int. Conf. Learn. Represent.*, 2024. 2, 5, 6, 3
- [9] Fadri Furrer, Tonci Novkovic, Marius Fehr, Abel Gawel, Margarita Grinvald, Torsten Sattler, Roland Siegwart, and Juan Nieto. Incremental object database: Building 3d models from multiple partial observations. In *IEEE/RSJ Int. Conf. on Intell. Robots and Syst.*, pages 6835–6842. IEEE, 2018. 3, 5
- [10] Golnaz Ghiasi, Xiuye Gu, Yin Cui, and Tsung-Yi Lin. Scaling open-vocabulary image segmentation with image-level labels. In *Eur. Conf. Comput. Vis.*, pages 540–557. Springer, 2022. 2
- [11] Margarita Grinvald, Fadri Furrer, Tonci Novkovic, Jen Jen Chung, Cesar Cadena, Roland Siegwart, and Juan Nieto. Volumetric instance-aware semantic mapping and 3d object discovery. *IEEE Robotics and Automation Letters*, 4(3): 3037–3044, 2019. 1, 2, 3, 5
- [12] Jun Guo, Xiaojian Ma, Yue Fan, Huaping Liu, and Qing Li. Semantic gaussians: Open-vocabulary scene understanding with 3d gaussian splatting. *arXiv preprint arXiv:2403.15624*, 2024. 2, 5, 6
- [13] Huy Ha and Shuran Song. Semantic abstraction: Open-world 3d scene understanding from 2d vision-language models. In *Conf. Robot. Learn.*, pages 643–653, 2022. 2

- [14] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Int. Conf. Comput. Vis.*, pages 2961–2969, 2017. 2
- [15] Chenguang Huang, Oier Mees, Andy Zeng, and Wolfram Burgard. Visual language maps for robot navigation. In *IEEE Int. Conf. on Robotics and Automation*, pages 10608–10615, 2023. 1
- [16] Rui Huang, Songyou Peng, Ayca Takmaz, Federico Tombari, Marc Pollefeys, Shiji Song, Gao Huang, and Francis Engelmann. Segment3d: Learning fine-grained class-agnostic 3d segmentation without manual labels. In *Eur. Conf. Comput. Vis.*, pages 278–295. Springer, 2024. 2, 5, 6
- [17] Krishna Murthy Jatavallabhula, Alihusein Kuwajerwala, Qiao Gu, Mohd Omaha, Tao Chen, Alaa Maalouf, Shuang Li, Ganesh Iyer, Soroush Saryazdi, Nikhil Keetha, et al. Conceptfusion: Open-set multimodal 3d mapping. In *Robotics: Science and Systems*, 2023. 1, 2, 3
- [18] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139–1, 2023. 2
- [19] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lerf: Language embedded radiance fields. In *Int. Conf. Comput. Vis.*, pages 19729–19739, 2023. 2
- [20] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Int. Conf. Comput. Vis.*, pages 4015–4026, 2023. 1, 2
- [21] Sebastián Barbas Laina, Simon Boche, Sotiris Papatheodorou, Simon Schaefer, Jaehyung Jung, and Stefan Leutenegger. Findanything: Open-vocabulary and object-centric mapping for robot exploration in any environment. *arXiv preprint arXiv:2504.08603*, 2025. 1
- [22] Boyi Li, Kilian Q Weinberger, Serge Belongie, Vladlen Koltun, and René Ranftl. Language-driven semantic segmentation. In *Int. Conf. Learn. Represent.*, 2022. 2
- [23] Feng Li, Hao Zhang, Peize Sun, Xueyan Zou, Shilong Liu, Chunyuan Li, Jianwei Yang, Lei Zhang, and Jianfeng Gao. Segment and recognize anything at any granularity. In *Eur. Conf. Comput. Vis.*, pages 467–484. Springer, 2024. 2
- [24] Wei Li, Junhua Gu, Benwen Chen, and Jungong Han. Incremental instance-oriented 3d semantic mapping via rgb-d cameras for unknown indoor scene. *Discrete Dynamics in Nature and Society*, 2020(1):2528954, 2020. 2
- [25] Yue Li, Qi Ma, Runyi Yang, Huapeng Li, Mengjiao Ma, Bin Ren, Nikola Popovic, Nicu Sebe, Ender Konukoglu, Theo Gevers, et al. Scenesplat: Gaussian splatting-based scene understanding with vision-language pretraining. In *Int. Conf. Comput. Vis.*, 2025. 2
- [26] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yinan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted clip. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 7061–7070, 2023. 2

- 596 [27] Tomas Berriel Martins, Martin R Oswald, and Javier Civera.
597 Ovo-slam: Open-vocabulary online simultaneous localization
598 and mapping. *CoRR*, 2024. 2, 5, 6
- 599 [28] Yang Miao, Iro Armeni, Marc Pollefeys, and Daniel Barath.
600 Volumetric semantically consistent 3d panoptic mapping.
601 In *IEEE/RSJ Int. Conf. on Intell. Robots and Syst.*, pages
602 12924–12931. IEEE, 2024. 1, 2, 4, 5
- 603 [29] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik,
604 Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf:
605 Representing scenes as neural radiance fields for view syn-
606 thesis. *Commun. ACM*, 65(1):99–106, 2021. 2
- 607 [30] Laksh Nanwani, Kumaraditya Gupta, Aditya Mathur,
608 Swayam Agrawal, AH Abdul Hafez, and K Madhava Kr-
609 ishna. Open-set 3d semantic instance maps for vision lan-
610 guage navigation–o3d-sim. *Advanced Robotics*, 38(19–20):
611 1378–1391, 2024. 2
- 612 [31] Richard A Newcombe, Shahram Izadi, Otmar Hilliges,
613 David Molyneaux, David Kim, Andrew J Davison, Pushmeet
614 Kohi, Jamie Shotton, Steve Hodges, and Andrew Fitzgib-
615 bon. Kinectfusion: Real-time dense surface mapping and
616 tracking. In *IEEE International Symposium on Mixed and
617 Augmented Reality*, pages 127–136. IEEE, 2011. 1, 4, 5
- 618 [32] Helen Oleynikova, Zachary Taylor, Marius Fehr, Roland
619 Siegwart, and Juan Nieto. Voxblox: Incremental 3d eu-
620 clidean signed distance fields for on-board mav planning. In
621 *IEEE/RSJ Int. Conf. on Intell. Robots and Syst.*, pages 1366–
622 1373. IEEE, 2017. 1
- 623 [33] Songyou Peng, Kyle Genova, Chiyu Jiang, Andrea
624 Tagliasacchi, Marc Pollefeys, Thomas Funkhouser, et al.
625 Openscene: 3d scene understanding with open vocabularies.
626 In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 815–824,
627 2023. 2, 5, 6, 7
- 628 [34] Jens Piekenbrinck, Christian Schmidt, Alexander Hermans,
629 Narunas Vaskevicius, Timm Linder, and Bastian Leibe.
630 Opensplat3d: Open-vocabulary 3d instance segmentation us-
631 ing gaussian splatting. In *IEEE Conf. Comput. Vis. Pattern
632 Recog.*, pages 5246–5255, 2025. 2
- 633 [35] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J
634 Guibas. Pointnet++: Deep hierarchical feature learning on
635 point sets in a metric space. *Adv. Neural Inform. Process.
636 Syst.*, 30, 2017. 2
- 637 [36] Lu Qi, Jason Kuen, Weidong Guo, Tiancheng Shen, Jiuxiang
638 Gu, Jiaya Jia, Zhe Lin, and Ming-Hsuan Yang. High-quality
639 entity segmentation. *arXiv preprint arXiv:2211.05776*, 2022.
640 3, 5, 8
- 641 [37] Minghan Qin, Wanhua Li, Jiawei Zhou, Haoqian Wang, and
642 Hanspeter Pfister. Langsplat: 3d language gaussian splatting.
643 In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 20051–
644 20060, 2024. 2
- 645 [38] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
646 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
647 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning
648 transferable visual models from natural language super-
649 vision. In *Int. Conf. Mach. Learn.*, pages 8748–8763. PMLR,
650 2021. 2
- 651 [39] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang
652 Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman
Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting
Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-
Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feicht-
enhofer. Sam 2: Segment anything in images and videos.
arXiv preprint arXiv:2408.00714, 2024. 2
- [40] Jonas Schult, Francis Engelmann, Alexander Hermans, Or
Litany, Siyu Tang, and Bastian Leibe. Mask3d: Mask trans-
former for 3d semantic instance segmentation. In *IEEE Int.
Conf. on Robotics and Automation*, pages 8216–8223, 2023.
2, 5
- [41] Nur Muhammad Mahi Shafiullah, Chris Paxton, Lerrel
Pinto, Soumith Chintala, and Arthur Szlam. Clip-fields:
Weakly supervised semantic fields for robotic memory.
arXiv preprint arXiv:2210.05663, 2022. 2
- [42] Julian Straub, Thomas Whelan, Lingni Ma, Yufan Chen, Erik
Wijmans, Simon Green, Jakob J Engel, Raul Mur-Artal, Carl
Ren, Shobhit Verma, et al. The replica dataset: A digital
replica of indoor spaces. *arXiv preprint arXiv:1906.05797*,
2019. 2, 5, 6, 3
- [43] Ayça Takmaz, Elisabetta Fedele, Robert W Sumner, Marc
Pollefeys, Federico Tombari, and Francis Engelmann. Open-
mask3d: Open-vocabulary 3d instance segmentation. In *Adv.
Neural Inform. Process. Syst.*, 2023. 2, 4, 5, 6, 7
- [44] Ayca Takmaz, Alexandros Delitzas, Robert W Sumner,
Francis Engelmann, Johanna Wald, and Federico Tombari.
Search3d: Hierarchical open-vocabulary 3d segmentation.
IEEE Robotics and Automation Letters, 2025. 1, 2
- [45] Silvan Weder, Hermann Blum, Francis Engelmann, and
Marc Pollefeys. Labelmaker: automatic semantic label gen-
eration from rgb-d trajectories. In *Int. Conf. on 3D Vis.*, pages
334–343. IEEE, 2024. 1
- [46] Luis Wiedmann, Luca Wiehe, and David Rozenberszki. Dc-
seg: Decoupled 3d open-set segmentation using gaussian
splatting. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages
5217–5226, 2025. 2
- [47] Yunyang Xiong, Bala Varadarajan, Lemeng Wu, Xiaoyu Xi-
ang, Fanyi Xiao, Chenchen Zhu, Xiaoliang Dai, Dilin Wang,
Fei Sun, Forrest Iandola, et al. Efficientsam: Leveraged
masked image pretraining for efficient segment anything.
In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 16111–
16121, 2024. 2
- [48] Xiaoxu Xu, Yitian Yuan, Jinlong Li, Qiudan Zhang, Ze-
qun Jie, Lin Ma, Hao Tang, Nicu Sebe, and Xu Wang. 3d
weakly supervised semantic segmentation with 2d vision-
language guidance. In *Eur. Conf. Comput. Vis.*, pages 87–
104. Springer, 2024. 2
- [49] Kashu Yamazaki, Taisei Hanyu, Khoa Vo, Thang Pham,
Minh Tran, Gianfranco Doretto, Anh Nguyen, and Ngan
Le. Open-fusion: Real-time open-vocabulary 3d mapping
and queryable scene representation. In *IEEE Int. Conf. on
Robotics and Automation*, pages 9411–9417. IEEE, 2024. 1,
2
- [50] Dianyi Yang, Yu Gao, Xihan Wang, Yufeng Yue, Yi Yang,
and Mengyin Fu. Opengs-slam: Open-set dense semantic
slam with 3d gaussian splatting for object-level scene under-
standing. *arXiv preprint arXiv:2503.01646*, 2025. 2
- [51] Yunhan Yang, Xiaoyang Wu, Tong He, Hengshuang Zhao,

- 710 and Xihui Liu. Sam3d: Segment anything in 3d scenes. *arXiv preprint arXiv:2306.03908*, 2023. [2](#)
- 711 [52] Mingqiao Ye, Martin Danelljan, Fisher Yu, and Lei Ke. Gaussian grouping: Segment and edit anything in 3d scenes. In *Eur. Conf. Comput. Vis.*, pages 162–179. Springer, 2024.
- 712 [53] Yingda Yin, Yuzheng Liu, Yang Xiao, Daniel Cohen-Or, Jingwei Huang, and Baoquan Chen. Sai3d: Segment any instance in 3d scenes. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3292–3302, 2024. [2](#)
- 713 [54] Xuan Yu, Yili Liu, Chenrui Han, Sitong Mao, Shunbo Zhou, Rong Xiong, Yiyi Liao, and Yue Wang. Panopticrecon: Leverage open-vocabulary instance segmentation for zero-shot panoptic reconstruction. In *IEEE/RSJ Int. Conf. on In-tell. Robots and Syst.*, pages 12947–12954. IEEE, 2024. [1](#), [2](#)
- 714 [55] Xuan Yu, Yuxuan Xie, Yili Liu, Haojian Lu, Rong Xiong, Yiyi Liao, and Yue Wang. Leverage cross-attention for end-to-end open-vocabulary panoptic reconstruction. *arXiv preprint arXiv:2501.01119*, 2025. [1](#)
- 715 [56] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Int. Conf. Comput. Vis.*, pages 11975–11986, 2023. [2](#), [5](#), [1](#)
- 716 [57] Jianhao Zheng, Daniel Barath, Marc Pollefeys, and Iro Armeni. Map-adapt: Real-time quality-adaptive semantic 3d maps. In *Eur. Conf. Comput. Vis.*, pages 220–237. Springer, 2024. [2](#), [4](#)
- 717 [58] Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and Andrew J Davison. In-place scene labelling and understanding with implicit scene representation. In *Int. Conf. Comput. Vis.*, pages 15838–15847, 2021. [2](#)
- 718 [59] Xueyan Zou, Zi-Yi Dou, Jianwei Yang, Zhe Gan, Linjie Li, Chunyuan Li, Xiyang Dai, Harkirat Behl, Jianfeng Wang, Lu Yuan, et al. Generalized decoding for pixel, image, and language. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 15116–15127, 2023. [2](#)
- 719 [60] Xueyan Zou, Jianwei Yang, Hao Zhang, Feng Li, Linjie Li, Jianfeng Wang, Lijuan Wang, Jianfeng Gao, and Yong Jae Lee. Segment everything everywhere all at once. *Adv. Neural Inform. Process. Syst.*, 36:19769–19782, 2023. [2](#)
- 720
- 721
- 722
- 723
- 724
- 725
- 726
- 727
- 728
- 729
- 730
- 731
- 732
- 733
- 734
- 735
- 736
- 737
- 738
- 739
- 740
- 741
- 742
- 743
- 744
- 745
- 746
- 747
- 748
- 749

CVI-MAP: Open-Vocabulary Instance-Semantic Mapping

Supplementary Material

6. Additional Details of the Method

6.1. Super-Point Merging

In the main paper, we addressed under-segmentation in 2D instances. Over-segmentation, however, is handled by merging spatially proximate super-points. Here, we provide details on the merging procedure.

For all voxels containing 3D points in the point cloud $P_{t,j}$, let $\Omega_{j,k}$ denote the number of voxels assigned to the super-point S_k with label k . The spatial proximity $Spa(S_a, S_b)$ measures the overlap between two super-points S_a and S_b by counting how many point clouds significantly intersect both:

$$Spa(S_a, S_b) = \sum_{P_j \in \mathcal{P}} \mathbb{I}[\Omega_{j,a} > \theta_{\text{assoc}} \wedge \Omega_{j,b} > \theta_{\text{assoc}}],$$

where $\mathcal{P} = \{P_{t,j}\}_{j=1}^{N_t}, \forall t \in T$ is the set of all inserted point clouds, and $\mathbb{I}[\cdot]$ is the Iverson bracket which is one if the condition inside holds and zero otherwise. Super-points with overlap $Spa(S_a, S_b)$ exceeding θ_{merge} are considered spatially connected, and merged into a single instance with same label.

6.2. TSDF Map Projection via Ray-Casting

After the incremental TSDF fusion and label stabilization with super-points voting and merging, we obtain a global TSDF-based instance map. We leverage this global instance map with stabilized instance labels across frames for the view-selection and feature extraction.

We obtain the projection of the TSDF map within the current camera frame by casting ray going through each pixel to the TSDF voxels. Combining with the depth prior from the depth input for ray-casting, we get a globally aligned 2D instance mask that taking occlusion and multi-view consistency into account.

7. Ablation Study on the VLM Backbone

Table 7 compares several VLM backbones for semantic feature extraction. We observe that larger SigLIP variants [56] yield consistent improvements in mIoU and AP metrics, with moderate parameter growth. Our method uses siglip-large-patch16-384, which offers the best trade-off between accuracy and computational cost.

8. Runtime

Table 8 reports the per-frame runtime of the main components, measured on an Nvidia RTX3090 GPU and an In-

Method	mIoU	mAcc	AP ₂₅	AP ₅₀	AP _{all}	Param.
clip-vit-large-patch14-336	15.4	24.7	20.6	14.9	6.7	0.4B
siglip-large-patch16-384 (Ours)	26.9	33.2	36.4	22.0	8.6	0.7B
siglip-so400m-patch14-384	25.2	33.8	32.4	19.1	9.7	0.9B
siglip2-large-patch16-384	<u>26.7</u>	<u>33.9</u>	<u>35.9</u>	<u>22.2</u>	10.4	0.9B
siglip2-so400m-patch14-384	26.4	<u>34.3</u>	<u>36.4</u>	<u>23.1</u>	9.4	1B

Table 7. **Semantic feature extraction** by (1) ViT-L, (2) siglip-large (Ours), (3) siglip-so400m, (4) siglip2-large, and (5) siglip2-so400m. The last column shows the number of parameters for each model.

Component	Run-Time / Frame (ms)	Thread
RGB Segmentation	964.8	1
Depth Segmentation	88.4	
2D-3D Association	76.3	2
View Selection	140.8	
Feature Extraction	131.3	3

Table 8. **Runtime breakdown.** All components run in parallel across dedicated threads. These steps are performed only on keyframes. This pipeline design keeps the overall system real-time despite potentially expensive individual modules.

tel Core i7-12700K CPU. The system operates as a multi-threaded pipeline, where segmentation, 2D-3D association, view selection, and semantic feature extraction run in parallel on separate CPU and GPU threads. The 2D-3D association includes 3D lifting of the 2D segments, update of the super-point map, and obtaining the global instance map via ray-casting.

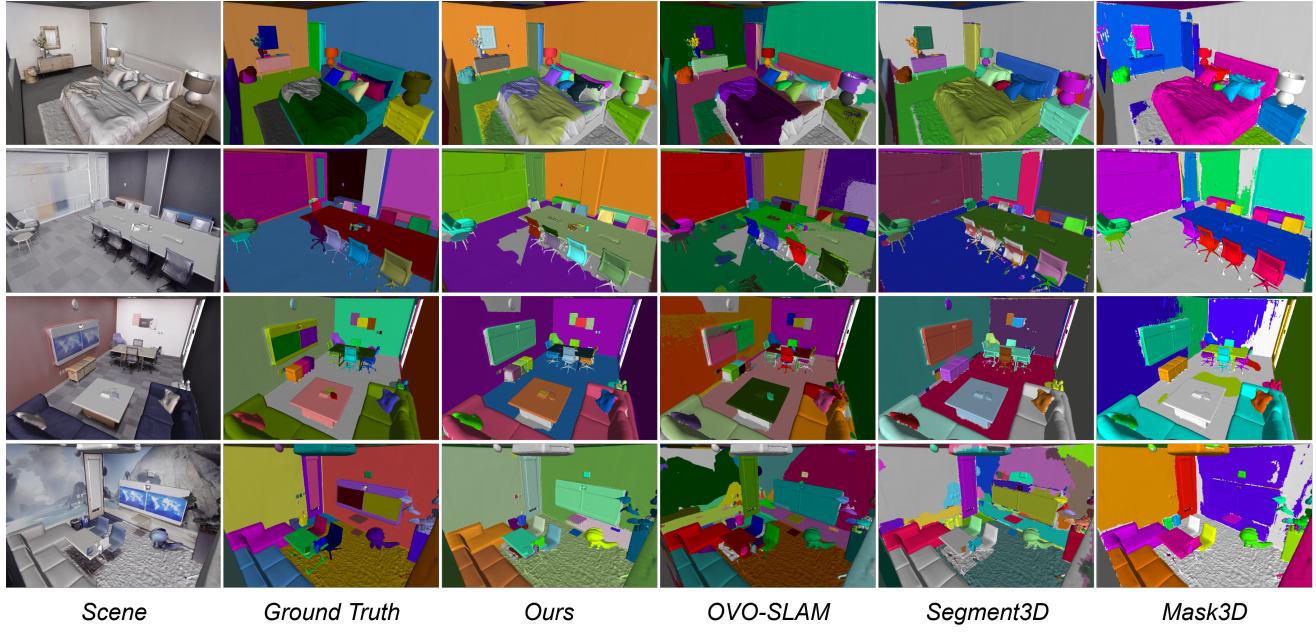
These instance/semantic segmentation steps are executed only on every n frames as reported in the main paper rather than every incoming frame, which substantially reduces overall latency.

9. Visualizations of Instance Maps

Figure 7 compares instance maps created by different methods. Online methods (Ours, OVO-SLAM) operate incrementally, while offline methods (Segment3D, Mask3D) process complete meshes. Colors correspond to instance IDs; gray denotes unobserved (for online methods) or unlabeled areas (for offline methods). Our method maintains consistent instance boundaries and achieves dense scene coverage, while offline methods leave large unlabeled regions despite full-scene access.

10. Heat Maps

To better understand how the objects are recognized, the heat maps of the instances in the scene are shown in Fig.8. It shows us how well can the objects we query can be distinguished from the others. We calculate the cosine similarities between the semantic features of all instances and the



(a) Qualitative comparison of instance maps on the Replica dataset.



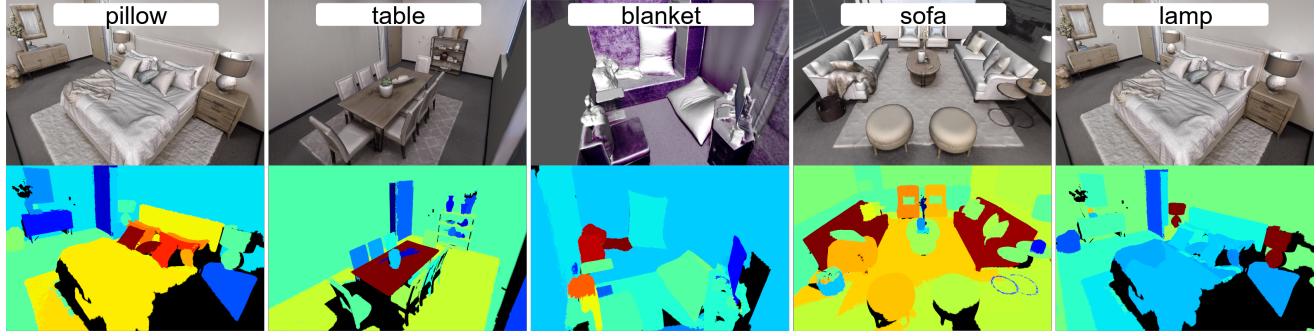
(b) Qualitative comparison of instance maps on the ScanNet dataset.

Figure 7. We compare our method with online [27] and offline [16, 40] approaches on the **Replica** (a) and **ScanNet** (b) datasets. Colors are randomly assigned for all instance maps according to the instance labels. Gray regions indicate unobserved areas for online methods (Ours and OVO-SLAM), and unlabeled for offline methods (Segment3D, Mask3D). OVI-MAP produces spatially coherent reconstructions, maintaining sharp instance boundaries throughout incremental mapping.

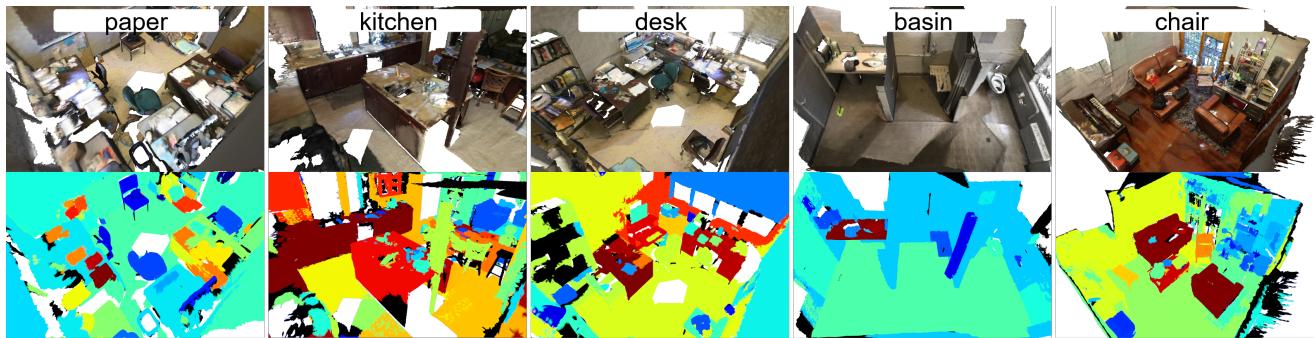
818 semantic features of the query semantic label, then map the
 819 normalized similarities to a color map.

820 The visualizations using heat maps further demonstrate
 821 the effectiveness of the proposed method in recognizing the
 822 objects in the scene, where the objects that match the query

823 label are highlighted correctly, and other irrelevant objects
 824 are not highlighted as much. This application shows the po-
 825 tential of the proposed method in real-world scenarios. For
 826 example, the system can locate the objects we are looking
 827 for based on the heat map, and update the semantic map



(a) Heat maps for semantic querying to the scenes from the Replica dataset.



(b) Heat maps for semantic querying to the scenes from the ScanNet dataset.

Figure 8. Heat map visualizations of the semantic queries. The color closer to red indicates the instance is more similar to the query semantic label, while the color closer to blue indicates the instance is less similar to the query semantic label. Unobserved areas are shown in black.

828

accordingly in those areas.

829

11. Datasets

830
831
832

We perform experiments on the Replica dataset [42] and ScanNet [6], which are widely used benchmarks for 3D scene understanding tasks.

833
834
835
836

We use 8 scenes including '*office0*', '*office1*', '*office2*', '*office3*', '*office4*', '*room0*', '*room1*', '*room2*' for evaluation on Replica, similar to OpenNeRF [8]. We evaluated on 18 scenes for ScanNet, as selected by ConceptFusion [17].