# Do Visual-Language Maps Capture Latent Semantics?

Matti Pekkanen, Tsvetomila Mihaylova, Francesco Verdoja, and Ville Kyrki

*Abstract*— **Visual-language models (VLMs) have recently been introduced in robotic mapping by using the latent representations, *i.e.*, embeddings, of the VLMs to represent the natural language semantics in the map. The main benefit is moving beyond a small set of human-created labels toward open-vocabulary scene understanding. While there is anecdotal evidence that maps built this way support downstream tasks, such as navigation, rigorous analysis of the quality of the maps using these embeddings is lacking. We investigate two critical properties of map quality: queryability and consistency. The evaluation of queryability addresses the ability to retrieve information from the embeddings. We investigate two aspects of consistency: intra-map consistency and inter-map consistency. Intra-map consistency captures the ability of the embeddings to represent abstract semantic classes, and inter-map consistency captures the generalization properties of the representation. In this paper, we propose a way to analyze the quality of maps created using VLMs, which forms an open-source benchmark to be used when proposing new open-vocabulary map representations. We demonstrate the benchmark by evaluating the maps created by two state-of-the-art methods, VLMaps and OpenScene, using two encoders, LSeg and OpenSeg, using real-world data from the Matterport3D data set. We find that OpenScene outperforms VLMaps with both encoders, and LSeg outperforms OpenSeg with both methods.**

## I. INTRODUCTION

Mobile robots must understand the geometry and semantics of the environment to accomplish complex tasks. The ability of maps to capture semantics has increased with the development of computer vision systems, from detecting single objects to creating pixel-wise semantic segmentation of images, yielding dense semantic maps with labels in each map cell [1]. However, most semantic segmentation methods are trained to segment the image into a small pre-selected set of categories. This imposes challenges in real-world settings when a robot operates in an environment that the categories do not fully describe [2], [3].

Visual-Language Models (VLMs) are networks that jointly train a visual and language encoder to learn a mapping from text and image inputs into a common visual-language latent space. Modern, efficient transformer architectures [4], such as CLIP [5], are trained with large data sets, which enables these models to effectively have an *open vocabulary*, meaning visual semantics can be matched to natural language instead of a set of selected symbols. Therefore, the

M. Pekkanen, T. Mihaylova, F. Verdoja and V. Kyrki are with School of Electrical Engineering, Aalto University, Espoo, Finland. {firstname.lastname}@aalto.fi
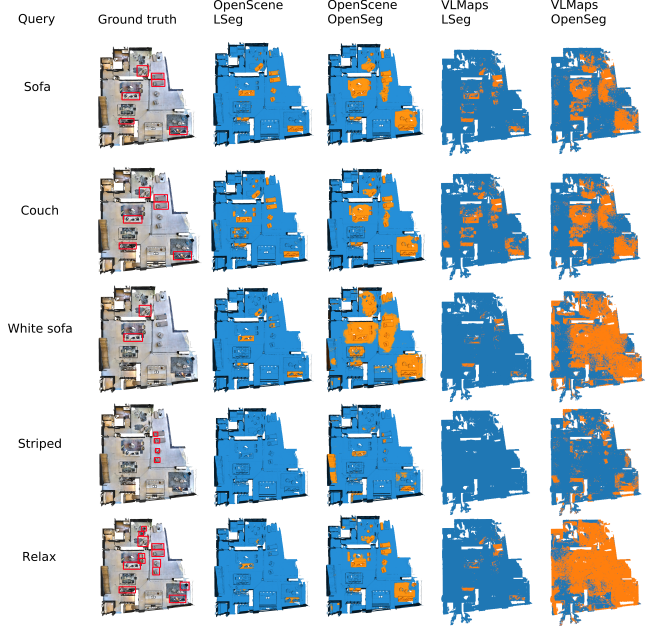
Fig. 1: Maps built from visual-language models have the ability to represent complex semantics, encompassing both the class of objects and their properties. We propose a benchmark for these type of maps and evaluate the quality of different state-of-the-art methods according to the consistency of their representation and their ability to be queried.

exploitation of VLMs in robotic mapping could allow richer semantics to be represented and to overcome the difficulty of adapting to new environments. While this work focuses on maps that encode semantics using VLMs, mapping methods have been proposed that can incorporate other modalities as well, such as audio [6]. We will collectively refer to all these maps as *latent semantic maps*.

The semantic information stored in the map must be able to be *queried*. While complex queries have been long possible in semantic maps by inference using object ontologies [7], semantic labels support only binary comparison as they are enumerations with no notion of distance. The promise of latent semantic maps is that queries can be evaluated by directly comparing their embeddings to the embeddings of the map with a similarity metric, yielding virtually unrestricted query vocabulary. Several successful algorithms have been built on VLMs that can perform mobile robot tasks, such as navigation [8]–[10]. However, the analysis from the robotic mapping perspective is left anecdotal.

In this work, we propose a way to evaluate the quality of latent semantic maps. We consider that three aspects of

the map capture its quality: queryability, consistency, and the breadth of the vocabulary. We concentrate on evaluating queryability and consistency because, to the best of our knowledge, no data sets exist to evaluate robotic maps that adequately capture open vocabulary semantics. Furthermore, we evaluate the consistency of the representation, *i.e.*, the embeddings themselves.

Understanding the possibilities and limitations of the representation is essential for bringing the methods into real-world applications. Having comparable metrics for the quality of the latent semantic maps is vital for measuring the progress in this new research field, as this allows for evaluating new methods against the state-of-the-art. For this reason, we provide an open-source benchmark[1] against which any new open vocabulary semantic representation can be tested. To demonstrate our benchmark, we evaluate two state-of-the-art methods: VLMaps [8] and OpenScene [11], using the Matterport3D [12] data set.

The main contributions of this paper are:

i) We define queryability and consistency of representation as measures of mapping quality for latent semantic maps and propose metrics to quantitatively evaluate these properties;

ii) We provide an open-source benchmark software for latent semantic maps, enabling evaluation of the progress in this new field of research;

iii) We provide a comprehensive analysis of two state-of-the-art methods using the proposed benchmark.

## II. RELATED WORK

### A. Maps created with visual-language embeddings

Creating maps using embeddings of VLMs, especially CLIP [5], has recently gained attention, with the proposal of methods such as VLMaps [8], OpenScene [11], NLMaps [13], and Uni-Fusion [14].

VLMaps and OpenScene create a dense grid map, where each cell is associated with a single VLM embedding. The embeddings are created from RGBD images with a language-driven semantic segmentation encoder, such as LSeg [15] or OpenSeg [16], both based on CLIP, and back-projected to the 3D map. Each cell is represented as the mean of the embeddings projected to the cell. In the original work, the 3D map in VLMaps is projected to 2D by averaging the embeddings in the $z$-axis. Additionally, OpenScene learns an encoder that directly produces CLIP embeddings from the 3D point cloud. Combining these two approaches, they propose joint 2D-3D "ensemble" features as the final representations in the map.

Unlike VLMaps and OpenScene, NLMaps creates a feature-based map, where each feature is associated with an VLM embedding. Uni-Fusion proposes a general method to create continuous maps from any perceptual data modality, including VLM embeddings. In this paper, we focus on

VLMaps and OpenScene because the NLMaps and Uni-Fusion source code were not made fully publicly available at the time of writing.

### B. Open-vocabulary queries

Open-vocabulary methods aim to move beyond predicting a constrained, *i.e.*, closed, set of labels. This means that not only should the open-vocabulary systems be able to capture synonyms and closely related terms but also attributes and related actions of the objects. This has been explored early on by formalizing the relation of labels into conceptual graphs [17], similar to what is used in semantic mapping, predicting attributes of objects [18] to the prediction of affordances and activities related to objects [11].

There are closely related fields to open-vocabulary detection and segmentation, especially open-set and zero-shot learning. In open-set learning, the task is to classify the known classes seen in training and to classify unseen classes as unseen [3]. Zero-shot learning aims to be able to classify classes that did not appear in the training data [19]. This, however, does not imply that open vocabulary is in use. Additionally, many open-vocabulary methods [20]–[22], even when trained with open-vocabulary methods, are evaluated in a closed-vocabulary but zero-shot setting. While this fails to capture the full potential of the open vocabulary, the results are comparable to semantic segmentation methods, where currently, fully supervised methods outperform the visual-language methods [11].

In contrast, in this work, we evaluate queryability such that the evaluation is not restricted to semantic segmentation of the map, which imposes that the map must be partitioned. Instead, each query is evaluated separately.

### C. Analysis of maps

When assessing the quality of robotic maps, most existing works concentrate on evaluating geometry using methods ranging from simple feature descriptors [23], frequency-based methods [24], to probabilistic measures [25].

Evaluating the geometric and semantic properties jointly is only relevant when the geometry and semantics are jointly estimated, *e.g.*, in semantic Simulatenous Localization and Mapping (SLAM). Often, the metrics are qualitative, but quantitative metrics such as the localization accuracy or map reconstruction quality compared to ground truth objects are used as well [26]–[28]. Like many semantic mapping works, [29]–[31], we focus on the evaluation of the semantic properties separately, using binary classification metrics.

In VLMaps, as with other zero-shot navigation methods [9], [10], the map quality is not directly assessed, but instead, they use the success of the downstream navigation task as the metric, which does not fully evaluate the capabilities of the representation.

Because the metrics used in evaluating latent semantic maps in each prior work are different, comparing the quality of their maps is non-trivial, which is the problem we address in this paper.

---

[1]The source code will be made publicly available upon acceptance of this work.

## III. PROBLEM STATEMENT

We aim to create a benchmark to evaluate the quality of latent semantic maps. We only consider methods where the estimation of geometry and semantics are disjoint, so we concentrate on evaluating the semantic quality of the maps. We believe that three aspects of the map capture its semantic quality: queryability, consistency, and the breadth of the vocabulary. We concentrate on evaluating queryability and consistency because, to the best of our knowledge, no data sets exist to evaluate robotic maps that adequately capture open-vocabulary semantics.

The benchmark must evaluate *queryability*, as it is the primary way to retrieve information from the representation and, therefore, acts as the measure of two things: first, indirectly, the ability of the representation to contain relevant information, and second, the accessibility of that information. It also must address *consistency*, as consistent performance within and between measuring runs, times, environments, and sensors is desirable.

## IV. METHODS

### A. Evaluating queryability

Because the environment consists both of background regions (*e.g.*, floor and ceiling) and discrete object instances (*e.g.*, a chair or a table), we propose evaluating queryability in two ways: by measuring the voxel-based coverage of the query results and the instance retrieval capability.

*1) Voxel-based queryability:* The first method evaluates the overall matching between the query results compared to the ground truth in a binary classification setting. This property cannot be evaluated using standard multi-class classification, as queries do not form a partition of the map. Each query produces a query result, a binary mask, over the whole map, and the same voxel might be a match of multiple queries. This is evident from the open-vocabulary perspective; *e.g.*, a sofa might match queries "sofa", "couch", "place to sit", and "soft".

The method consists of the following steps:

i) A map $m$ is created from each sequence from the data set, forming the set $\mathcal{M}$.

ii) Each map $m \in \mathcal{M}$ is queried with each query $q$ in the set of queries $\mathcal{Q}$. The query result of a single query is a binary segmentation mask of the map, $\hat{\mathbf{y}}_q = \{\hat{y}_1, \ldots, \hat{y}_N\}$, consisting of $N$ voxel masks $\hat{y} \in \{true, false\}$. Combined, the query results form the set of voxel-based predictions $\hat{Y}_V = \{\hat{\mathbf{y}}_q \forall q \in \mathcal{Q}\}$.

iii) For each query, the true mask $\mathbf{y}_q$ is created from the ground truth map $m_g$ such, that the voxels answering query $q$ are $true$, others $false$. The true masks form a set of true labels $Y_V = \{\mathbf{y}_q \forall q \in \mathcal{Q}\}$.

iv) The binary classification metrics are calculated between $Y_V$ and $\hat{Y}_V$.

Each comparison of a prediction $\hat{y}$ to the corresponding ground truth $y$ yields either true positive $tp = \hat{y} \wedge y$, true negative $tn = \neg\hat{y} \wedge \neg y$, false positive $fp = \hat{y} \wedge \neg y$, or false negative $fn = \neg\hat{y} \wedge y$. Using these, binary classification

metrics are defined: accuracy $a = \frac{tp+tn}{tp+tn+fp+fn}$, precision $p = \frac{tp}{tp+fp}$, recall $r = \frac{tp}{tp+fn}$, F1-score $f_1 = \frac{2tp}{2tp+fp+fn}$, and Intersection over Union (IoU) $iou = \frac{tp}{tp+fn+tp}$.

As the regions of interest are relatively small with most of the queries, true negative predictions dominate the predictions. Therefore, accuracy cannot be used as a metric, as it incorporates true negative predictions. The other metrics are not affected by true negative predictions, so we use them as the metrics for binary classification tasks in this work.

*2) Instance-based queryability:* The second method evaluates the capability of the map to detect and retrieve all matching objects for the query. Each instance is predicted to match or not match the query; therefore, this measures the coverage of matches within an object rather than over the whole map.

While some latent semantic maps propose a way to perform instance segmentation, not all do, so measuring instance segmentation performance is outside of the scope of this work. Instead, we assume given ground-truth instances; these can be obtained through many different approaches, using latent semantic maps alone or combining them with other maps or modalities.

The method consists of the following steps:

i) Given $\mathcal{M}$, $\mathcal{Q}$, and ground truth instance segmentation $\mathcal{I}$, each map $m \in \mathcal{M}$ is queried with the each query $q \in \mathcal{Q}$, each query yielding binary predicted mask $\hat{\mathbf{y}}_q$.

ii) For each object instance $i \in \mathcal{I}$, the corresponding voxels $\hat{\mathbf{i}}_{i,q} \subset \hat{\mathbf{y}}_q$ are selected from the predicted map.

iii) If the majority of the voxels in $\hat{\mathbf{i}}_{i,q}$ are $true$, the prediction $\hat{y}_{i,q}$ is $true$, otherwise $false$. The predictions for each instance for each query combined form the set of predictions $\hat{Y}_I = \{\hat{y}_{i,q} \forall i \in \mathcal{I}; q \in \mathcal{Q}\}$.

iv) Similar to the previous method, the true map $m_q$ is created. If the instance $i$ on the true map answers the query $q$, the true label $y_{i,q} = true$, otherwise $false$. The true labels for each instance for each query form the set of all true labels denoted $Y_I$.

v) The binary classification metrics are calculated between $Y_I$ and $\hat{Y}_I$.

From $Y_I$ and $\hat{Y}_I$, the same metrics presented for the voxel-based queryability are computed for each instance.

### B. Evaluating consistency

We further subdivide the consistency of the map into two distinct aspects, intra-map consistency and inter-map consistency, and propose a method for evaluating each.

*1) Intra-map consistency:* Intra-map consistency captures the similarity of embeddings across the voxels within a map, sharing semantic meaning. The hypothesis is that embeddings sharing semantic meaning are clustered together in the latent space to allow the separability of different concepts. While each object instance is in some sense unique, a consistent set of embeddings represents an abstract base class to which the instances belong.

The method consists of the following steps:

i) Given a map $m$, and semantic label $l$, $\mathcal{E}_l$ is the set of embeddings corresponding to voxels in the map $m$

where the ground truth semantic label is $l$, and $\mathcal{E}_m$ is the set of all embeddings in the map $m$. We form a set of tuples $\mathcal{T} = \{(\mathcal{E}_l, \mathcal{E}_m) \forall m \in \mathcal{M}; l \in \mathcal{L}\}$, where $\mathcal{L}$ is a closed-set semantic label vocabulary. For computational efficiency and to no considerable change in results, in practice, the sets $\mathcal{E}_l$ and $\mathcal{E}_m$ are subsampled. We use a subsampling ratio of 0.1 in this work.

ii) The average absolute deviation $d_l^m$ is calculated for $\mathcal{E}_l$, and $d_m^m$ for $\mathcal{E}_m$, for each tuple $t \in \mathcal{T}$. The average absolute deviation $d_a$ measures the statistical dispersion of a set, and given a set of points $\mathcal{E} = \{e_1, ..., e_n\}$ it is defined as

$$d_a = \frac{1}{n} \sum_{j=1}^{n} |f_d(e_j, \bar{\mathcal{E}})|, \qquad (1)$$

where $f_d$ is a deviation metric and $\bar{\mathcal{E}}$ a central point of set $\mathcal{E}$. We use the mean of the embeddings as the central point and cosine distance as the deviation metric, which is consistent with the use of cosine distance loss for CLIP.

iii) The intra-map consistency ratio $c_l^m = \frac{d_l^m}{d_m^m}$ is computed for each tuple $t \in \mathcal{T}$, with label $l$ in map $m$. This ratio represents the distinguishability of the class compared to the map average.

*2) Inter-map consistency:* Inter-map consistency captures the similarity of voxels with the same semantic label across different maps. The hypothesis is that embeddings within the same label are closer to each other across maps than to those with different labels with respect to a distance metric. This would imply that the objects retain their distinctiveness across maps, and therefore, the system generalizes better across different environments.

The method consists of the following steps:

i) Given the set of tuples $\mathcal{T}$, constructed according to Section IV-B.1.

ii) For all pairs of tuples $((\mathcal{E}_{l,1}, \mathcal{E}_{m,1}), (\mathcal{E}_{l,2}, \mathcal{E}_{m,2}))$, $\mathcal{E}_{m,1} \neq \mathcal{E}_{m,2}$, parametric Wasserstein 2-distance $d_w$ is calculated between $\mathcal{E}_{l,1}$ and $\mathcal{E}_{l,2}$.

The Wasserstein $p$-distance is computationally heavy for large sets of high-dimensional embeddings. For this reason, given an n-dimensional embedding $e \in \mathbb{R}^n$, we approximate the set of embeddings $\mathcal{E}$ with an $n$-dimensional normal distribution $\mathcal{N}(\mu, P)$. This allows us to have a closed-form solution for the Wasserstein 2-distance

$$d_w(\mathcal{E}_1, \mathcal{E}_2) \approx \|\mu_1 - \mu_2\|_2^2 + \mathrm{tr}\left(P_1 + P_2 - 2(P_2^{\frac{1}{2}} P_1 P_2^{\frac{1}{2}})^{\frac{1}{2}}\right). \quad (2)$$

## V. EXPERIMENTS

The main motivation of the experiments is to demonstrate the proposed benchmark. The two main questions the benchmark aims to answer with the experiments are:

1) How queryable state-of-the-art latent semantic maps are?

2) How consistent are their visual-language embeddings within and across maps?

To answer these questions, we evaluated two state-of-the-art methods, VLMaps [8] and OpenScene [11], in a series of experiments.

### A. Data set

The data set used in the benchmark is the Matterport3D data set [12], which consists of houses captured with an RGBD camera and ground truth semantic segmentation. The data set was used in the original work of VLMaps, in which they selected a sample of 10 sequences from the data set. In the benchmark, we use the same sequences and the same set of 42 labels forming the set $\mathcal{L}$; their names are also used in the following experiments as the query set $\mathcal{Q}$ to leverage the semantic ground truth provided with the dataset[2].

Since Matterport3D does not provide ground truth instance segmentation, instances used as ground truth were segmented using a region growing algorithm based on the ground truth semantics [32]. Each voxel is initialized as a seed cluster; then, region growing steps are performed, where the labels of all neighboring clusters are compared. If the labels are the same, the clusters are joined. Otherwise, they are not. This step is iterated until no more clusters can be joined or a maximum iteration limit is reached.

### B. Parametrization of the methods

Instead of using the 2D mapping method proposed in the original work [8], the VLMaps map was created using the new 3D mapping method available at the repository of the original authors [33]. This choice was made to create a fair comparison using state-of-the-art methods. The central idea is the same, except the 3D map avoids the problems when aggregating the embeddings along the $z$-axis. All of the parameters were the default used by the original authors. The 3D grid voxel size used is 0.05 m, as proposed by the authors.

The 3D mapping method was also extended to create semantic ground truth maps from the Matterport3D data. The VLMaps use Habitat simulator [34]–[36] for camera measurements. Similarly to the RGB images, we create semantic images and backproject them to the 3D environment.

When queried with a query $q$, VLMaps creates a list of queries by setting $q$ and string "other" into a list of phrases such as "A photo of ..." and "there is ... in the scene". Then, the most similar query is selected from the list for each voxel. The binary mask is $true$ for voxels where the most similar query is a phrase containing $q$; otherwise, it is $false$. The binary mask is then dilated to encompass whole object instances.

OpenScene maps were created using the parameters proposed in the original paper. We use the proposed 2D-3D ensemble features. The 3D grid voxel size used is 0.02 m, as proposed by the authors. OpenScene does not provide a way to query open-vocabulary binary masks. Because we use the set of labels as the queries, the binary query mask

---

[2]Supplementary material including the list of labels will be available in arXiv and link provided here upon acceptance of this work.

TABLE I: Average results of the voxel- and instance-based queryability tests over 10 maps.

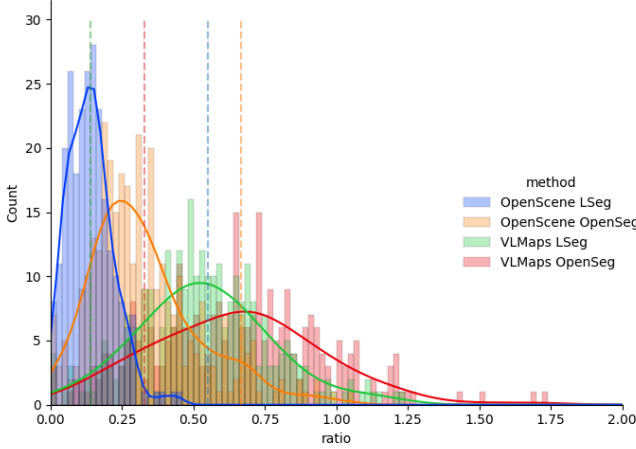| Method | Encoder | Voxel-based | | | | Instance-based | | | |
|--------|---------|-------------|-----------|--------|-------|----------------|-----------|--------|-------|
| | | F1 | Precision | Recall | IoU | F1 | Precision | Recall | IoU |
| OpenScene | LSeg | **0.623** | **0.617** | 0.630 | **0.457** | 0.258 | 0.252 | 0.264 | 0.149 |
| OpenScene | OpenSeg | 0.580 | 0.575 | 0.587 | 0.412 | **0.261** | **0.255** | 0.267 | **0.150** |
| VLMaps | LSeg | 0.498 | 0.401 | **0.661** | 0.332 | 0.257 | 0.205 | **0.347** | 0.147 |
| VLMaps | OpenSeg | 0.393 | 0.287 | 0.625 | 0.246 | 0.202 | 0.144 | 0.342 | 0.112 |



Fig. 2: The intra-map consistency ratios of the methods as histograms, with a kernel density estimation, which is smoothed with a Gaussian kernel. The mean of the distribution is depicted with a vertical line of corresponding color.

is created by comparing the equality of the voxel labels and the query.

Both methods used the same pre-trained LSeg and OpenSeg encoder provided by the original authors. Both encoders are based on the CLIP backbone; LSeg uses the CLIP-ViT-B/32 backbone with 512-dimensional embeddings, whereas OpenSeg uses the CLIP-ViT-L/14 with 768-dimensional embeddings.

*C. Results*

*1) Queryability:* The results of the queryability benchmark are presented in Table I, where the F1-score, precision, recall, and IoU of the method and encoder combinations are presented for the voxel- and instance-based tests. Overall, OpenScene performs better than VLMaps with both encoders, based on the higher F1-scores and IoUs. The higher recall shown by VLMaps with LSeg is likely the result of the post-processing of the queries, where the query results are dilated to encompass whole objects. The instance-based results mostly align with the voxel-based results, except for the slightly decreased performance in all metrics.

While OpenScene with LSeg is better in the voxel-based experiment, both encoders have virtually the same performance in the instance-based experiment when used within OpenScene. This aligns with the findings presented in [11], where OpenScene with LSeg encoder was found to have

better IoU.

Therefore, these results show that both the mapping method and choice of encoder matter in the map creation process. They also indicate that the 3D structure of the environment, included by OpenScene, contains information that can be leveraged in addition to the purely image-based creation of embeddings.

All methods have better recall than precision, which are usually tradeoffs. Better recall might be beneficial for queryability, especially in a navigation context: having more object candidates to explore than missing potential objects is better. This is especially visible in VLMaps, where the recall is considerably higher than precision. OpenScene has a more balanced approach, reflected in the high F1-Score, which is a good metric for the overall performance of the methods. IoU also conforms to the results, increasing their reliability.

As the instance classification is based on the prediction of the majority label, many objects seem to match the query only with a region representing less than half of the object. In this regard, the query methodology and, subsequently, the experiment setup could be improved: instead of classifying the results with binary masks, the minimum cosine distance could be used. This could allow the detection of instances with a small but significant match to the query. This would be relevant, for instance, when the object is partially occluded or unobserved. For example, if a sofa is partially covered with a blanket, the image-based methods could segment it as a blanket, but with the visible small region, it could still be recognized as a sofa. In this regard, the 3D understanding could prove extremely valuable: The geometry of the sofa might not significantly change when it is covered, so it can still be recognized as a sofa. When using open vocabulary, the description could even cover both objects: a sofa covered with a blanket.

Additionally, we qualitatively demonstrate the capability of the benchmark to process open-vocabulary queries. We use the VLMaps' binary masking method without prompt engineering or dilation as the comparison method. From the results presented in Figure 1, where the matches are shown in orange and non-matches in blue, it can be seen that OpenScene query results are considerably less noisy and encompass entire objects, indicating that the features are more distinctive. While LSeg yields the objects more accurately, the more abstract queries, such as "striped" or "relax", produce no results. OpenScene seems to be more sensitive, finding solutions to abstract queries but yielding many false positive results. While this is in line with the

(a) OpenScene with LSeg

(b) OpenScene with OpenSeg

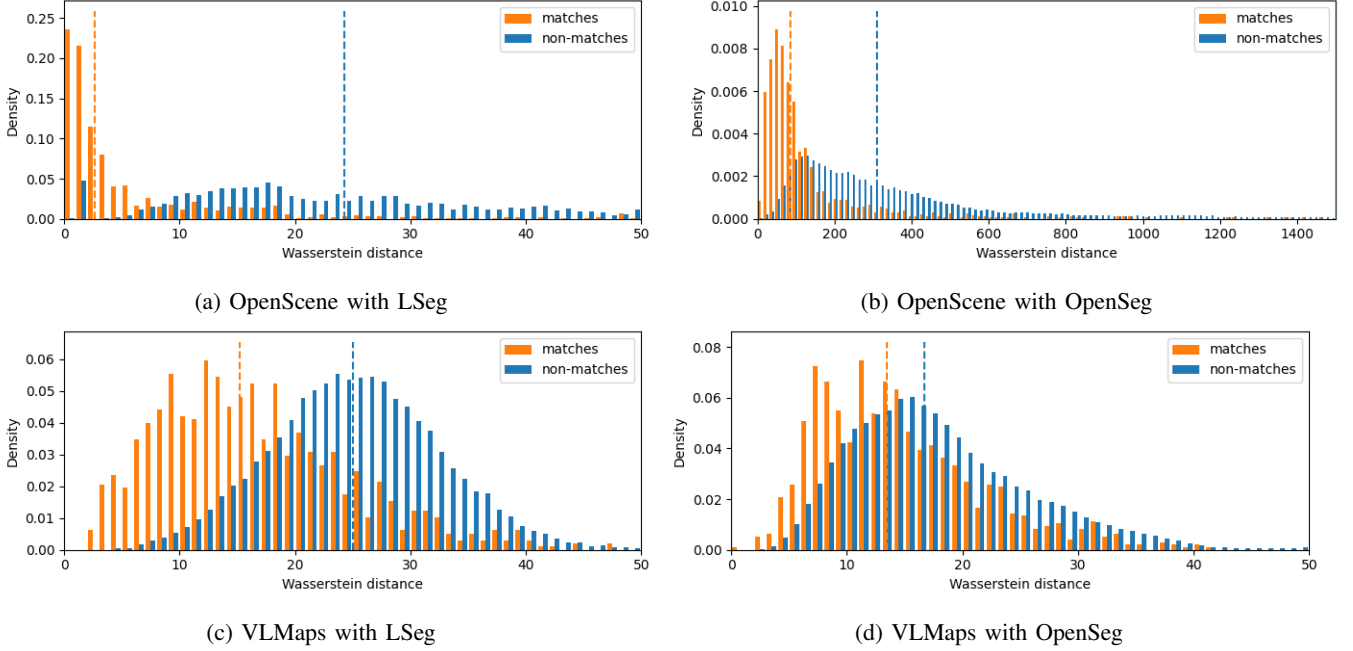(c) VLMaps with LSeg

(d) VLMaps with OpenSeg

Fig. 3: The distribution of Wasserstein distances of matching and non-matching labels are shown in orange and blue, respectively. The median of each distribution is presented with a dashed line.

findings of the other experimental results, it might hint towards OpenSeg having a broader vocabulary.

*2) Consistency:*

*a) The intra-map consistency:* The results are presented in Figure 2, in which the intra-map consistency ratios of the methods are illustrated with histograms. Additionally, the means of data are shown with dashed lines. Notably, the ordering of the methods according to their intra-map consistency is the same as in the queryability experiment presented in Section V-C.1. OpenScene has better consistency between the labels with each encoder, which is indicated by the lower means of the distributions. This indicates that OpenScene can cluster the embeddings of labels better using the ensemble features. VLMaps with OpenSeg has a long tail that continues beyond the figure, with several outliers.

However, there are multiple classes with a ratio above one. This means that the embeddings in a class are more different from each other than the map on average. This means these classes are challenging to cluster and hard to distinguish. However, these classes have, on average, 2-4% the number of voxels compared to the average number of voxels in the whole data set. For example, the two most represented categories in the tail are "misc" and "objects", both relatively non-descriptive generic categories, and therefore, this could be attributed to the limited number of classes in the experiment.
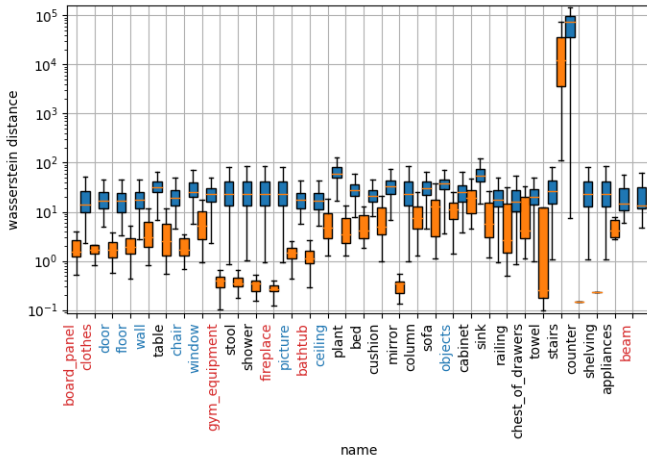
*b) The inter-map consistency:* In Figure 3, the Wasserstein distances between matching, *i.e.*, the label is the same, and non-matching sets of embeddings across all maps are shown. The medians of the data are shown with dashed lines. The ratio between the medians of matching and non-matching labels suggests that OpenScene separates the labels

better than VLMaps, and LSeg separates them better than OpenSeg. The ratio of medians of OpenScene with LSeg is 9.21, OpenScene with OpenSeg is 3.67, VLMaps with LSeg 1.65, and VLMaps with OpenSeg 1.24, which is once again the exact ordering observed in Sections V-C.1 and V-C.2.a.
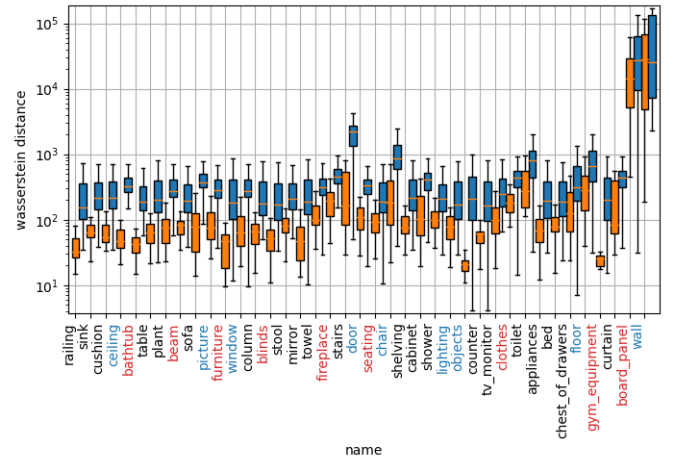
While the shapes of distributions are similar, OpenScene has much larger distances on average, even with matching labels. This results from the fact that the distilled features of OpenSeg are of a larger magnitude than the image features of both VLMaps and OpenScene image features. For example, the average norm of all OpenScene with OpenSeg embeddings in the map of the first sequence is 13.68 times larger than the average norm of all VLMaps with OpenSeg embeddings in the same map, and the maximum being 57.97 times larger than the maximum VLMaps feature.

Furthermore, in Figure 4, the distances of matching and non-matching distances are shown per class with orange and blue bars, respectively. The same phenomena can be seen: LSeg distinguishes the classes better than OpenSeg and OpenScene better than VLMaps, which is indicated by the separation of the blue and orange box plots. Additionally, the ten most common and uncommon categories in the Matterport3D data set are marked with red and blue names, respectively. VLMaps performs well in separating common classes and worse with uncommon classes, especially using LSeg, whereas OpenScene performs better across the whole set of labels.
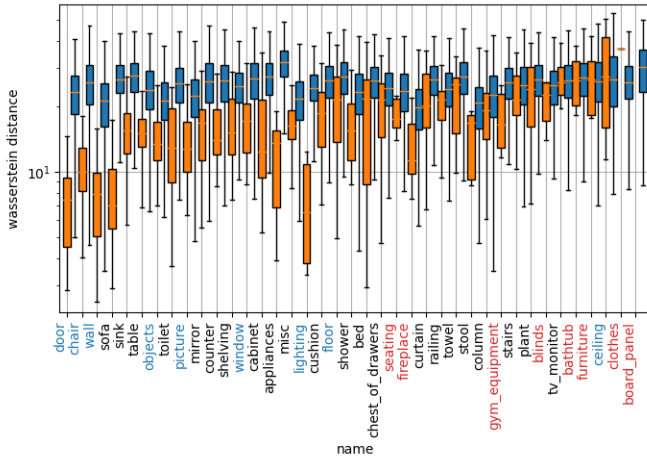
When considering all these results, the combination of OpenScene and LSeg produces maps that are both more queryable and consistent, making it the method that currently can be expected to generalize better and lead to more consistent performance in downstream robotic applications.
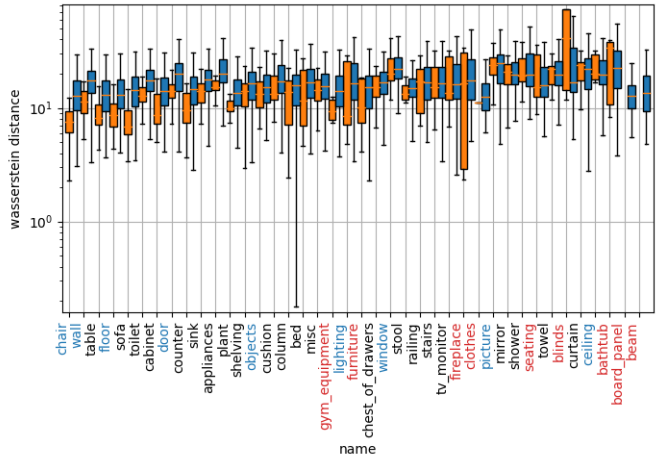
(a) OpenScene LSeg



(b) OpenScene OpenSeg



(c) VLMaps LSeg



(d) VLMaps OpenSeg

Fig. 4: The Wasserstein distances of the matching labels are depicted as orange bars and non-matching labels as blue bars on a logarithmic scale. The ten most common and uncommon labels are marked with red and blue text, respectively. The labels are sorted per image on the separability of the distance distribution according to Krusal-Wallis statistic [37]. The outliers of the boxplots have been omitted for clarity. The OpenScene distances, in (a) and (b), are multiple orders of magnitude larger, as the embeddings are an order of magnitude larger.

## VI. CONCLUSION

Mobile robot maps must be able to capture the semantics of the environment for the robots to perform complex tasks. If the semantic representation is constrained to a closed set of classes, the robot's ability to understand the environment is constrained to the same set. Latent semantic maps present a way to overcome this limitation and enable open-vocabulary scene understanding by representing the semantics in a latent space with a notion of distance.

In this work, we propose quantitative metrics to evaluate the quality of state-of-the-art latent semantic maps. In our experiments, OpenScene mapping using 3D scene structure with LSeg semantic embedding performed the best. This indicates that the 3D structure of the environment is useful for inferring the latent semantics. While LSeg performed best in the semantic segmentation setting, the qualitative demonstration hints that OpenSeg might be more sensitive to rarer and more abstract queries when properly thresholded.

The ability of visual-language maps to capture latent semantics is, however, still limited. Currently, fully supervised semantic segmentation methods outperform all latent semantic maps. The relatively low F1-scores and IoUs, especially in the instance-based test, indicate that there is still work to be done to address the problems of thresholding and clustering the query results to segment object instances.

A data set providing ground truth open-vocabulary semantics is direly needed. This would allow extending the evaluation presented here to adequately address the breadth of vocabulary and capture the full promise of VLMs.

## REFERENCES

[1] S. Garg *et al.*, "Semantics for Robotic Mapping, Perception and Interaction: A Survey," *Foundations and Trends® in Robotics*, vol. 8, no. 1–2, pp. 1–224, 2020, arXiv:2101.00443 [cs].

[2] A. Bendale and T. E. Boult, "Towards Open Set Deep Networks," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Las Vegas, NV, USA: IEEE, Jun. 2016, pp. 1563–1572.

[3] C. Geng, S.-J. Huang, and S. Chen, "Recent Advances in Open Set Recognition: A Survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 43, no. 10, pp. 3614–3631, Oct. 2021.

[4] A. Vaswani *et al.*, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[5] A. Radford *et al.*, "Learning transferable visual models from natural language supervision," in *International Conference on Machine Learning*. PMLR, 2021, pp. 8748–8763.

[6] C. Huang, O. Mees, A. Zeng, and W. Burgard, "Audio visual language maps for robot navigation," 2023, arXiv:2303.07522 [cs].

[7] M. Tenorth and M. Beetz, "KNOWROB - knowledge processing for autonomous personal robots," in *2009 IEEE/RSJ International Conference on Intelligent Robots and Systems*. St. Louis, MO, USA: IEEE, oct 2009, pp. 4261–4266.

[8] C. Huang, O. Mees, A. Zeng, and W. Burgard, "Visual Language Maps for Robot Navigation," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. London, United Kingdom: IEEE, may 2023, pp. 10 608–10 615.

[9] S. Y. Gadre, M. Wortsman, G. Ilharco, L. Schmidt, and S. Song, "Cows on pasture: Baselines and benchmarks for language-driven zero-shot object navigation," in *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Vancouver, Canada, jun 2023, pp. 23 171–23 181.

[10] D. Shah, B. Osiński, S. Levine *et al.*, "Lm-nav: Robotic navigation with large pre-trained models of language, vision, and action," in *Conference on Robot Learning*. Atlanta, GA, USA: PMLR, nov 2023, pp. 492–504.

[11] S. Peng *et al.*, "Openscene: 3d scene understanding with open vocabularies," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Vancouver, Canada, jun 2023, pp. 815–824.

[12] A. Chang *et al.*, "Matterport3D: Learning from RGB-D Data in Indoor Environments," Sep. 2017, arXiv:1709.06158 [cs].

[13] B. Chen *et al.*, "Open-vocabulary queryable scene representations for real world planning," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. London, UK: IEEE, may 2023, pp. 11 509–11 522.

[14] Y. Yuan and A. Nüchter, "Uni-fusion: Universal continuous mapping," *IEEE Transactions on Robotics*, vol. 40, pp. 1373–1392, jan 2024.

[15] B. Li, K. Q. Weinberger, S. Belongie, V. Koltun, and R. Ranftl, "Language-driven Semantic Segmentation," Apr. 2022, arXiv:2201.03546 [cs].

[16] G. Ghiasi, X. Gu, Y. Cui, and T.-Y. Lin, "Scaling open-vocabulary image segmentation with image-level labels," in *European Conference on Computer Vision*. Tel Aviv, Israel: Springer, oct 2022, pp. 540–557.

[17] H. Zhao, X. Puig, B. Zhou, S. Fidler, and A. Torralba, "Open Vocabulary Scene Parsing," in *2017 IEEE International Conference on Computer Vision (ICCV)*. Venice: IEEE, Oct. 2017, pp. 2021–2029.

[18] M. A. Bravo, S. Mittal, S. Ging, and T. Brox, "Open-vocabulary Attribute Detection," in *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Vancouver, BC, Canada: IEEE, Jun. 2023, pp. 7041–7050.

[19] A. Bansal, K. Sikka, G. Sharma, R. Chellappa, and A. Divakaran, "Zero-Shot Object Detection," in *Computer Vision – ECCV 2018*, V. Ferrari, M. Hebert, C. Sminchisescu, and Y. Weiss, Eds., vol. 11205. Cham: Springer International Publishing, sep 2018, pp. 397–414, series Title: Lecture Notes in Computer Science.

[20] A. Zareian, K. D. Rosa, D. H. Hu, and S.-F. Chang, "Open-Vocabulary Object Detection Using Captions," in *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Nashville, TN, USA: IEEE, Jun. 2021, pp. 14 388–14 397.

[21] X. Gu, T.-Y. Lin, W. Kuo, and Y. Cui, "Open-vocabulary Object Detection via Vision and Language Knowledge Distillation," May 2022, arXiv:2104.13921 [cs].

[22] C. Feng *et al.*, "Promptdet: Towards open-vocabulary detection using uncurated images," in *European Conference on Computer Vision*. Tel Aviv, Israel: Springer, oct 2022, pp. 701–717.

[23] A. I. Wagan, A. Godil, and X. Li, "Map quality assessment," in *Proceedings of the 8th Workshop on Performance Metrics for Intelligent Systems*. Gaithersburg Maryland: ACM, Aug. 2008, pp. 278–282.

[24] T. P. Kucner, M. Luperto, S. Lowry, M. Magnusson, and A. J. Lilienthal, "Robust Frequency-Based Structure Extraction," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. Xi'an, China: IEEE, May 2021, pp. 1715–1721.

[25] S. Aravecchia, M. Clausel, and C. Pradalier, "Comparing metrics for evaluating 3D map quality in natural environments," *Robotics and Autonomous Systems*, vol. 173, p. 104617, Mar. 2024.

[26] R. F. Salas-Moreno, R. A. Newcombe, H. Strasdat, P. H. Kelly, and A. J. Davison, "SLAM++: Simultaneous Localisation and Mapping at the Level of Objects," in *2013 IEEE Conference on Computer Vision and Pattern Recognition*. Portland, OR, USA: IEEE, Jun. 2013, pp. 1352–1359.

[27] J. Mccormac, R. Clark, M. Bloesch, A. Davison, and S. Leutenegger, "Fusion++: Volumetric Object-Level SLAM," in *2018 International Conference on 3D Vision (3DV)*. Verona: IEEE, Sep. 2018, pp. 32–41.

[28] M. Runz, M. Buffier, and L. Agapito, "MaskFusion: Real-Time Recognition, Tracking and Reconstruction of Multiple Moving Objects," in *2018 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. Munich, Germany: IEEE, Oct. 2018, pp. 10–20.

[29] N. Sunderhauf, T. T. Pham, Y. Latif, M. Milford, and I. Reid, "Meaningful maps with object-oriented semantic mapping," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. Vancouver, BC: IEEE, Sep. 2017, pp. 5079–5085.

[30] W. Chen, S. Hu, R. Talak, and L. Carlone, "Leveraging Large (Visual) Language Models for Robot 3D Scene Understanding," Nov. 2023, arXiv:2209.05629 [cs].

[31] G. Narita, T. Seno, T. Ishikawa, and Y. Kaji, "PanopticFusion: Online Volumetric Semantic Mapping at the Level of Stuff and Things," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. Macau, China: IEEE, Nov. 2019, pp. 4205–4212.

[32] R. Adams and L. Bischof, "Seeded region growing," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 16, no. 6, pp. 641–647, Jun. 1994.

[33] Chenguang Huang and Oier Mees and Andy Zeng and Wolfram Burgard. Vlmaps. [Online]. Available: https://github.com/vlmaps/vlmaps

[34] M. Savva *et al.*, "Habitat: A platform for embodied ai research," in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, Seoul, Korea, oct 2019, pp. 9338–9346.

[35] A. Szot *et al.*, "Habitat 2.0: Training home assistants to rearrange their habitat," in *Advances in Neural Information Processing Systems*, M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, Eds., vol. 34. Curran Associates, Inc., 2021, pp. 251–266.

[36] X. Puig *et al.*, "Habitat 3.0: A co-habitat for humans, avatars and robots," 2023, arXiv:2310.13724 [cs].

[37] W. H. Kruskal and W. A. Wallis, "Use of ranks in one-criterion variance analysis," *Journal of the American statistical Association*, vol. 47, no. 260, pp. 583–621, apr 1952.