# Polyline Generative Navigable Space Segmentation for Autonomous Visual Navigation

Zheng Chen[1], Zhengming Ding[2], David J. Crandall[1], Lantao Liu[1]

*Abstract*—Detecting navigable space is a fundamental capability for mobile robots navigating in unknown or unmapped environments. In this work, we treat visual navigable space segmentation as a scene decomposition problem and propose Polyline Segmentation Variational autoencoder Network (PSV-Net), a representation learning-based framework for learning the navigable space segmentation in a self-supervised manner. Current segmentation techniques heavily rely on fully-supervised learning strategies which demand a large amount of pixel-level annotated images. In this work, we propose a framework leveraging a Variational AutoEncoder (VAE) and an AutoEncoder (AE) to learn a polyline representation that compactly outlines the desired navigable space boundary. Through extensive experiments, we validate that the proposed PSV-Net can learn the visual navigable space with no or few labels, producing an accuracy comparable to fully-supervised state-of-the-art methods that use all available labels. In addition, we show that integrating the proposed navigable space segmentation model with a visual planner can achieve efficient mapless navigation in real environments.

*Index Terms*—Label-efficient learning, variational autoencoders, segmentation, visual navigation.

## I. INTRODUCTION

**F**OR mobile robots to be able to navigate unknown spaces, it is crucial to understand the traversability of complex environments that consist of cluttered objects. The goal is to construct collision-free traversable space, which we term as *navigable space*. If cameras are used to perceive the environment, a typical way to identify navigable space is through image segmentation by leveraging deep neural networks to perform multi-class [1] or binary-class [2] segmentation of images. The present work takes a binary-class segmentation approach in which the robot needs to distinguish navigable space from non-navigable space in the robot's perceived first-person-view (FPV) images.

However, most existing deep neural network-based methods are developed on top of a fully-supervised learning paradigm and rely on annotated datasets such as Cityscapes [3] or

Fig. 1: Navigation in a real cluttered environment with the proposed PSV-Net. The ground robot is circled in the left image. The right image shows a motion primitive-based visual planner [5] that is used for navigation. The yellow arrow indicates the goal, the green paths are motion primitives, the red path is the optimal one, and the yellow path is the historically traversed path. SEDF represents the Scaled Euclidean Distance Field; see [5].

KITTI [4]. These datasets usually contain an immense number of pixel-level annotated segmented images. Collecting and annotating such data for robotic applications in novel environments is prohibitively costly and time-consuming.

To overcome the limitation of fully-supervised learning and pave a path for mobile robot navigation, we propose a self-supervised learning method by treating binary navigable space segmentation as a *scene decomposition* problem, and the training signals come from certain visual inputs such as surface normal images. It has been demonstrated in [6], [7] that scene decomposition can be solved in an unsupervised fashion with Variational AutoEncoders (VAEs). In contrast to semantic segmentation by supervised learning where the model is trained by human-annotated pixel-wise labels, scene decomposition attempts to learn the compositional nature from visual information alone without explicit annotations.

Most (if not all) existing VAE-based scene decomposition and representation learning methods use pixel-wise learning in which the value of every pixel is predicted. However, pixel-wise methods prone to generate segmentation with noises, e.g., pixel islands or scatters, which can inevitably affect downstream decision-making in the planning module, leading to the whole navigation system being vulnerable to unsafe/inefficient decisions. In addition, it is also expensive to adopt a post-processing module to smooth out possible pixel noises. An intuitive way to solve issues of pixel-wise methods is to directly learn the desired boundary for downstream planning tasks. The boundary could be represented by polylines/splines (represented by vertices or control points) [9], [10], which are naturally robust to noise in images.

The proposed PSV-Net consists of two networks, Net-I (a VAE) and Net-II (an AE). During training, the goal of Net-I is to learn pseudo labels from surface normals by using

Fig. 2: Framework overview. $I_S$, $E_1/E_2$, and $D_1/D_2$ are the input surface normal image, encoders, and decoders for both nets. **Net-I**: $q_\phi(z|x^{(i)})$ represents the latent categorical distribution from which a segmentation sample $z$ is sampled using ♠ – the Gumbel-Softmax sampler for discrete distributions proposed in [8]. Prior I is to regularize the "shape" of the predicted latent distribution. $R_1^{(i)}$ is the reconstructed input and $\mathcal{L}_1$ is the reconstruction loss between $I_S$ and $R_1^{(i)}$. **Net-II**: $\mathcal{V}$ are the initial vertices for polyline representation. $\mathcal{P}$ is a set of predicted vertices from $E_2$. $R_2^{(i)}$ is the reconstructed image and $\mathcal{L}_2$ is the reconstruction between $R_2^{(i)}$ and the latent image $I_q$, converted from the latent distribution $q_\phi$.

categorical distributions as the latent representation. Using the supervision signals from Net-I, Net-II learns a set of vertices which describe the location of the navigable space boundary. We show that our boundary-based PSV-Net achieves comparable segmentation results to fully-supervised learning-based state-of-the-art methods. To demonstrate the efficacy of our PSV-Net on real navigation tasks, we also integrate the proposed segmentation model with a visual planner to guide the robot to move without collision (see Fig. 1).

## II. RELATED WORK

Traditional image segmentation methods with no supervision focus on crafting features and energy functions to define desired objectives. For example, active contour-based models [11] optimize over a polygon (represented by vertices) by means of energy minimization based on both the image features and shape priors such as boundary continuity and smoothness. However, active contours lack flexibility and heavily rely on low-level image features and global parameterization of priors [12]. Recently, deep active contour-based models have been proposed [12], but these methods require ground-truth contour vertices and thus belong to the supervised learning paradigm. Another line of research uses adversarial approaches for unsupervised segmentation, e.g., work in [13] explores the idea of changing the textures or colors of objects without changing the overall distribution of the dataset, and proposes an adversarial architecture to segment a foreground object in each image. Although the adversarial methods show impressive results, they suffer from instabilities in training. The proposed technique in this paper is close in spirit to some work in scene decomposition [6] where the goal is to decompose scenes into objects in an unsupervised fashion with a generative model.

Our proposed framework is built upon VAEs, which are often used for representation learning [14], [15] to learn useful representations of data with little or no supervision. The trained generative model (decoder) can generate new images from any random sample in the learned latent distribution. Many variants of VAEs have been proposed following [16], where the latent representation is described as a continuous distribution. In this paper, to learn the navigable space, instead, we use a discrete distribution as the latent representation.

In addition, the value of boundaries in image segmentation has been shown in a large amount of previous literature. The polygon/spline representation is used in [10], [17]–[20] to achieve fast and potentially interactive instance segmentation. Acuna *et al* propose a new approach to learn sharper and more accurate semantic boundaries [21]. By treating boundary detection as the dual task of semantic segmentation, a new loss function with a boundary consistency constraint to improve boundary pixel accuracy for semantic segmentation is designed [22]. The work in [23] proposes a content-adaptive down-sampling technique that learns to favor sampling locations near semantic boundaries of target classes. By doing so, segmentation accuracy and computational efficiency can be balanced. Although the above methods make use of boundaries to improve performance, they all require ground truth labels for training — an important limitation that we aim to overcome in this work.

## III. METHODOLOGY

The proposed PSV-Net consists of two sub-nets, as shown in Fig. 2. Our intuition for this two-net architecture is based on a coarse-to-fine process, where Net-I learns a rough pixel-wise pseudo label of navigable space by taking advantage of generative modeling, while Net-II generates a refined version of the prediction from Net-I using a set of vertices.

### A. Net-I

The goal of Net-I is to learn a pixel-wise pseudo label from the input surface normal image. A standard VAE aims to learn a generative model which can synthesize new data from a random sample in a specified latent space. We want to obtain the parameters of the generative model by maximizing the data marginal likelihood $\log p_\theta(\mathbf{x}^{(i)})$, where $\mathbf{x}^{(i)}$ is a data point in our training dataset $\left\{\mathbf{x}^{(i)}\right\}_{i=1}^N$. Using Bayes' rule, the likelihood can be written,

$$\log p_\theta(\mathbf{x}^{(i)}) = \log \sum_{\mathbf{z}} \left[ p_\theta(\mathbf{x}^{(i)}|\mathbf{z}) p_\theta(\mathbf{z}) \right], \tag{1}$$

where $p_\theta(\mathbf{z})$ is the prior distribution of the latent representation and $p_\theta(\mathbf{x}^{(i)}|\mathbf{z})$ is the generative probability distribution of the reconstructed input given the latent representation. We can use

Fig. 3: Structure of Net-II. *I#1*: Surface normal image, *I#2*: Initial vertices, *I#3*: Newly predicted vertices, and *I#4*: Rendered output. $E_2$: encoder of Net-II. IFE: Image Feature Extraction net. GCN: Graph Convolutional Network. $D_2$: decoder of Net-II. NR: Neural Rendering module. In both blocks, blue points are vertices of polylines while orange points are auxiliary points on image bottom boundary for the convenience of neural rendering.



Fig. 4: GCN module in the *left block* of Fig. 3.

a neural network to approximate $p_\theta(\mathbf{x}^{(i)}|\mathbf{z})$ where $\theta$ can be thought as parameters of the network, but we cannot perform the sum operation over $\mathbf{z}$, hence Eq. (1) is computationally intractable.

An alternative is to find the lower bound of Eq. (1); to do this, we can use an encoder to approximate the true posterior of the latent variable $p_\theta(\mathbf{z}|\mathbf{x}^{(i)})$. We denote the encoder (another neural network) as $q_\phi(\mathbf{z}|\mathbf{x}^{(i)})$, where $\phi$ are parameters of the encoder network. Then we can derive the lower bound of Eq. (1) as [16],

$$\log p_\theta(\mathbf{x}^{(i)}) = \mathbb{E}_\mathbf{z}\left[\log p_\theta(\mathbf{x}^{(i)}|\mathbf{z})\right] - \mathbb{KL}(q_\phi(\mathbf{z}|\mathbf{x}^{(i)})||p_\theta(\mathbf{z})) +$$
$$\mathbb{KL}(q_\phi(\mathbf{z}|\mathbf{x}^{(i)})||p_\theta(\mathbf{z}|\mathbf{x}^{(i)})). \quad (2)$$

To perform segmentation, we assume $p_\theta(\mathbf{x}^{(i)}|\mathbf{z}) = \prod_{j=1}^{J}\mathcal{N}(\mathbf{x}^{(i)};R_1^{(i)},\sigma^2)$, $q_\phi(\mathbf{z}|\mathbf{x}^{(i)}) = \prod_{j=1}^{J}Cat(z_j|\mathbf{x}^{(i)},\phi)$, and the prior $p_\theta(\mathbf{z}) = \prod_{j=1}^{J}Cat(z_j)$, where $J$ is the number of pixels in the input image (same meaning in the following equations). Then only the last term $\mathbb{KL}(q_\phi(\mathbf{z}|\mathbf{x}^{(i)})||p_\theta(\mathbf{z}|\mathbf{x}^{(i)}))$ of Eq. (2) is unknown. Fortunately, we know that $\mathbb{KL}(p||q) \geq 0$ is always true for any two distributions $p$ and $q$. Therefore,

$$\log p_\theta(\mathbf{x}^{(i)}) \geq \mathbb{E}_\mathbf{z}\left[\log p_\theta(\mathbf{x}^{(i)}|\mathbf{z})\right] - \mathbb{KL}(q_\phi(\mathbf{z}|\mathbf{x}^{(i)})||p_\theta(\mathbf{z})), \quad (3)$$

where the term on the right side of the inequality is called the *Evidence Lower BOund* (ELBO) of the data marginal likelihood. Maximizing the ELBO is equivalent to maximizing the data marginal likelihood. By using Monte Carlo sampling, the total loss to be *minimized* over the whole training dataset is:

$$\mathcal{L}_I = \sum_{i=1}^{N}\left[\mathbb{KL}(q_\phi(\mathbf{z}|\mathbf{x}^{(i)})||p_\theta(\mathbf{z})) + \frac{1}{2KJ\sigma^2}\sum_{k=1}^{K}\left(\left\|\mathbf{x}^{(i)} - R_{1k}^{(i)}\right\|_2^2 + \frac{J}{2}\log\sigma^2\right)\right], \quad (4)$$

where $J$ is the total number of pixels, $N$ is the number of images in our training dataset, and $K$ is the number of samples in Monte Carlo sampling. $R_{1k}^i$ is the reconstructed input from the $k^{th}$ sampled $z$ in Monte Carlo sampling. The reconstruction loss $\mathcal{L}_1$ in Fig. 2 is equivalent to the term $\frac{1}{2K\sigma^2}\sum_{k=1}^{K}\left\|\mathbf{x}^{(i)} - R_1^{(i)}\right\|_2^2$.

For Net-I, we use a U-Net [24] as the $E_1$ to approximate the latent true posterior, and one convolutional layer as the $D_1$ to reconstruct an image from the latent code. The output of $E_1$, $q_\phi(\mathbf{z}|\mathbf{x}^{(i)})$, represents the categorical distribution, from which we can obtain a pesudo label by computing $I_q^{(i)} = \arg\max_c q_\phi(\mathbf{z}|\mathbf{x}^{(i)})$, where $c$ is the channel number of the output in $E_1$. We use $I_q^{(i)}$ as the pseudo label to supervise the training of Net-II, which we now describe.

### B. Net-II

With the supervision signals from Net-I, Net-II aims to learn a set of vertices which describe the location of the navigable space boundary. Net-II is an autoencoder, where the encoder manipulates the coordinates of a set of vertices (from $\mathcal{V}$ to $\mathcal{P}$ in Fig. 2). Then the predicted vertices $\mathcal{P}$ are used to render a reconstructed image $R_2^{(i)}$ by the decoder $D_2$. The training of the Net-II is supervised by the pseudo label $I_q^{(i)}$ from Net-I. The overall objective of Net-II, $\mathcal{L}_{II}$, is equal to the reconstruction loss $\mathcal{L}_2$ (see Fig. 2), which consists of a mean squared error and an appearance matching error [25] estimated using the Structural Similarity Index Measure (SSIM),

$$\mathcal{L}_{II} = \mathcal{L}_2 = \lambda_1 \frac{1 - SSIM(I_q^{(i)}, R_2^{(i)})}{2} + \lambda_2 \frac{1}{J}\left\|I_q^{(i)} - R_2^{(i)}\right\|_2^2, \quad (5)$$

where $J$ is the same as in Eq. (4), $\lambda_1$ and $\lambda_2$ are weights, and $\lambda_1 + \lambda_2 = 1$.

We use a hybrid net consisting of a convolutional neural network and a graph convolutional network as the encoder (see the *left block*, $E_2$ in Fig. 3), and a neural rendering module [26] as the decoder (see the *right block*, $D_2$ in Fig. 3). Specific Net-II components (shown in Fig. 3) are described as follows.

**Image Feature Extraction (IFE):** The IFE module is adopted to generate feature maps at different resolutions, where different layers have different resolutions. The feature maps are then used as features for vertices in the GCN module (will be discussed next), where we use the coordinates of the vertices to index locations of the features. Bilinear interpolation is applied to obtain the coordinates in different layers. The features from the IFE module give the downstream network more information regarding the vertices, such that the vertices' shift prediction can be well supported by visual clues.

**Graph Convolution Network (GCN):** We construct a graph with each node as the concatenation of the extracted

Fig. 5: Triangle selection for neural rendering. (a) Auxiliary points (orange), example predicted vertices (blue), and constructed triangles (black lines, only those formed by boundary points are considered in the selection step). (b) Orange lines connect predicted vertices with corresponding projected vertices on the image bottom border. Only those triangles which do not intersect with any of orange lines are ignored (e.g., the red one) to render. (c) Desired navigable space from neural rendering is marked as green region.



Fig. 6: Integration of the visual planner [5] with segmentation model. *Left:* Green curves: motion primitives in map space, Red curve: the optimal primitive, and Yellow arrow: Goal. *Middle:* Segmentation images at two time steps during navigation. *Right:* Scaled Euclidean Distance Field based on the segmentation images. The 2D projections of motion primitives are also visualized.

image features and their coordinates. Our GCN structure (see Fig. 4) is inspired by the network structure proposed in [19], [27]. The fundamental layer of the GCN module is the graph convolutional layer. We denote a graph as $\mathcal{G} = \{\mathcal{V}, \mathcal{E}, \mathcal{F}\}$, where $\mathcal{V} = \{v_i\}_{i=1}^{T}$ denotes the node set, $\mathcal{E} = \{e_j\}_{j=1}^{M}$ represents the edge set, and $\mathcal{F} = \{f_i\}_{i=1}^{T}$ is the feature vector set for nodes in the graph. $T$ and $M$ are the number of nodes and edges, respectively. Then a graph convolution layer is defined as,

$$f_i^{l+1} = \Gamma(w_0^l f_i^l + \sum_{j \in Nei(i)} w_1^l f_j^l), \qquad (6)$$

where $f_i^{l+1}$ and $f_i^l$ are the feature vectors on vertex $i$ before and after the convolution, $Nei(i)$ are the neighboring nodes of $i$, and $\Gamma(\cdot)$ is the activation function.

Note that the vertices are iteratively manipulated twice by the GCN, but the IFE only performs a single forward pass. In each manipulation, we use the updated coordinates to extract new features from the same set of feature maps.

**Neural Rendering (NR):** To convert the vertices predicted from the encoder $E_2$ to a *segmentation image* such that Net-II can be trained using the pseudo label from Net-I, we triangularize those vertices and select proper triangles for neural rendering. We first select three auxiliary vertices $\mathcal{V}'$ on the bottom boundary of image (see orange points of the *left* block of Fig. 3). The auxiliary vertices $\mathcal{V}'$ are fixed through the GCN module such that the input to $D_2$ is a new set of vertices consisting of $\mathcal{V}'$ and $\mathcal{P}$: $\mathcal{P}+ = \mathcal{P} \cup \mathcal{V}'$. The motivation for the auxiliary points is that we need to construct a *closed* area for neural rendering processing. We then use Delaunay triangulation to construct triangles on $\mathcal{P}+$. Next we use neural rendering [26] to render the constructed triangles into a segmentation output while keeping the whole pipeline differentiable. However, a discrepancy exists: Delaunay triangulation always returns a series of triangles forming a *convex* hull but the shape of real navigable space is not necessarily always convex, so the NR module tends to generate more triangles than needed. We propose to use a simple triangles selection method (see Fig. 5) to filter out unnecessary triangles.

### C. Visual Planning

To make the proposed model effective in real navigation, we integrate the PSV-Net with an image-based planner [5], which projects a library of motion primitives from map space

(see the left image of Fig. 6) to image space. This allows us to only use the binary segmentation images to make the decision on which primitive is optimal to track, leading to a *mapless* navigation system.

Specifically, first we convert the binary segmentation to a Scaled Euclidean Distance Field (SEDF); see two examples in right images of Fig. 6. Values in the SEDF represent costs of colliding with obstacle boundaries. Next, we compute a library of motion primitives [28], [29] $\mathcal{M} = \{\mathbf{p}_1, \mathbf{p}_2, \cdots, \mathbf{p}_n\}$ projected from map space to SEDF. Then we compute the navigation cost function for each primitive based on the evaluation on the SEDF and target progress (a distance measure from robot's current pose to the goal pose). Finally we select the primitive with the minimal cost to execute. The trajectory selection problem can be defined as,

$$\mathbf{p}_{optimal} = \underset{\mathbf{p}}{\text{argmin}} \; w_1 \cdot C_c(\mathbf{p}) + w_2 \cdot C_t(\mathbf{p}), \qquad (7)$$

where $C_c(\mathbf{p}) = \sum_j^m c_c^j$ and $C_t(\mathbf{p}) = \sum_j^m c_t^j$ are the collision cost and target cost of one primitive $\mathbf{p}$, and $w_1$, $w_2$ are corresponding weights, respectively. Readers can find more details about the SEDF and how to compute target progress in [5].

## IV. EXPERIMENTS

### A. Baselines and Evaluation Metrics

To evaluate the proposed method, we compare against two baselines, SNE-RoadSeg [2] and OFF-Net [30], both of which are pixel-wise, state-of-the-art methods for binary navigable space segmentation. We conduct extensive qualitative and quantitative experiments on multiple datasets. Note that our method takes depth/surface images as input. For fair comparison, all methods in our experiments require that at least depth images are provided. We use the mean of Intersection over Union (mIoU) as the metric for reporting quantitative results. The definition of IoU is: $\mathbf{IoU} = \frac{n_{tp}}{n_{tp} + n_{fp} + n_{fn}}$, where $n_{tp}, n_{tn}, n_{fp}$ and $n_{fn}$ are true positive, true negative, false positive, and false negative, respectively.

### B. Datasets

We evaluate the proposed method on several different datasets, including the standard KITTI road dataset [31], the

Fig. 7: Qualitative results on (a) KITTI and (b) Cityscapes. First row: RGB images; Second row: ground-truth navigable space segmentation. Labels for each row of prediction results (from the third row to the last row) have the form S/O/P-#, where S represents the baseline SNE-RoadSeg [2], O represents the baseline OFF-Net [30], P represents the proposed method PSV-Net, and # represents the percentage of the ground truth labels used. The baseline methods will not work if no ground truth labels are provided, so S-0 and O-0 are not available.

Cityscapes dataset [3], and the ORFD dataset for off-road environments [30].

The original KITTI road segmentation dataset [3] is specifically designed for binary navigable space segmentation in city-like environments. Since our proposed method learns segmentation from images in a self-supervised manner, we can only learn free space from images where geometric boundaries are consistent with human-defined labels. Otherwise, the self-learned segmentation will be different from the ground truth labels (but still reasonable), causing confusion for evaluation. For example, *sidewalk* sometimes shares the same geometric plane with *road*, meaning both classes can be geometrically navigable, but in human-predefined labels only the class *road* is marked as navigable. Therefore, throughout the experiments (including the baseline experiments), we only use the images starting with *um* and *umm* with corresponding *road* (instead of *lane*) ground truth labels in training data. We convert the LiDAR point clouds to depth images and use the surface normal estimation (SNE) module proposed in [2] to compute surface normal images.

Cityscapes [3] is a standard semantic segmentation dataset in which pixel-wise semantic labels are provided for images across multiple cities. We treat the *road* label as navigable while others as non-navigable, thus creating binary ground truth images. Note that we apply no data selection to the Cityscapes data as there is no way to distinguish images which have inconsistency between the geometry and the semantic label (of the *road*). This inevitably causes us to underestimate the performance of our proposed method. We use the provided

disparity images to compute depth images and then obtain normal images using SNE in [2].

ORFD [30] is a newly-released binary segmentation dataset for various off-road environments, e.g., farmland, grassland, woodland, countryside, and snowland. The data in ORFD have a notably different distribution than images in KITTI or Cityscapes, so we can test generalization capability during adaptation from one dataset to another, e.g., KITTI → ORFD (meaning the model is trained using KITTI but tested on ORFD).

To further validate our proposed method on various data, we also build a small dataset which is collected in indoor environments. We name the data we collect as INDOOR dataset. Examples of the indoor data can be seen in Fig. 9. Note that in INDOOR, only RGB and depth images are collected. The surface normal images are computed in the same way as [2]. The INDOOR dataset has no ground truth labels.

### C. Implementation Details

To implement PSV-Net as illustrated in Fig. 2, $E_1$ is a U-Net [24] with 2D convolution kernels and $D_1$ is a simple fully convolutional network. Similar to [27], IFE in $E_2$ is a VGG-16 [32] network up to layer *conv5_4* and we use the concatenation of features from layer *conv2_3, conv3_3, conv4_3* and *conv5_4* as the visual features of the nodes in the graph. The GCN in $E_2$ is based on [19], [27]. We use the pretrained mesh renderer proposed in [26] as $D_2$ (NR). We fix the parameters of $D_2$, and only parameters of $E_1$, $D_1$, and $E_2$ are updated during

training. We use the *Adam optimizer* for both networks and set the learning rate $lr_1 = 1e-3$ for Net-I and $lr_2 = 1e-4$ for Net-II. The total number of training epochs is 15 and we fix the weights of Net-I after epoch 3. We set the weights in $\mathcal{L}_2$ (Eq. (5)) as $\lambda_1 = 0.8$ and $\lambda_2 = 0.2$. We compare to the most recent baseline [2], a fully-supervised state-of-the-art method for predicting navigable space. To permit a fair comparison, we use the same batch size, 1, for both methods. We implement the whole framework using PyTorch [33] and conduct all experiments on a single Nvidia Geforce RTX 2080 Super GPU.

### D. Comparison for KITTI

A qualitative comparison of different methods on KITTI is shown in Fig. 7(a). Predictions are reported with 4 different percentages of ground truth labels ($0\%, 1\%, 30\%$, and $100\%$) used during training. We show that the proposed boundary based method is superior to the pixel-wise methods in Fig. 7(a). The baseline method SNE-RoadSeg produces a large number of noisy pixels when few labels are presented, e.g., results in the fourth row of Fig. 7(a). The baseline OFF-Net has a similar problem (see the fifth line) although the noise is much less than SNE-RoadSeg. This is because the training data is insufficient and the model is still underfitting to the data, leading to large prediction errors. However, pixel-wise methods can still show large amounts of noise even when we increase the number of labels in the training data; see the second- and third-to-the-last rows in Fig. 7(a). One could argue that the noise appears because the training data is insufficient, and that simply collecting more training data would fix the problem. However, in practice, it is very difficult to collect large amounts of training data, and even when large datasets are available, pixel-wise noise still easily occurs when we try to deploy a model in a new environment with a domain shift from the training data.

In contrast, we observe that the proposed boundary-based PSV-Net is able to yield estimates (the $3^{rd}$ row) close to the ground truth labels even without any ground truth labels provided during training. All models can achieve better predictions as the amount of ground truth increases, but PSV-Net still maintains an advantage over the baselines. There are few noisy pixels appearing in the resulting predictions, which in turn can significantly improve results of the downstream planning modules. A more detailed quantitative comparison can be seen in Fig. 8(a).

### E. Comparison for Cityscapes

We show a qualitative comparison for the Cityscapes dataset in Fig. 7(b), again comparing results across different amounts of available ground truth labels. We see that the navigable spaces in Cityscapes exhibit more complicated patterns than KITTI due to the existence of more objects in the scene (e.g., cars and pedestrians), but our proposed method can still detect the shape of the navigable space and outperform the baseline when no or few labels are available. A quantitative comparison is in Fig. 8(b). Note that the output of OFF-Net collapses with very few labels (1%), because the model is



Fig. 8: Quantitative comparisons of different methods as a function of the quantity of ground turth available during training, on (a) KITTI dataset and (b) Cityscapes dataset.

not able to learn a reasonable representation. Interestingly, the mIoU with 0% labels is greater than the mIoU with 1%. One possible explanation is that the 0% model generates random outputs based on randomly-initialized network weights, which are actually better than the model trained with 1%.

As the number of ground truth labels increases, the baseline results improve significantly while PSV-Net remains about the same. The reasons for this are two-fold. First, we do not apply data selection to Cityscapes, and there exist many images having inconsistency between geometry and semantics. While our method may generate geometrically-correct segmentation maps that are required for autonomous navigation, it may not match the semantics defined in the ground truth labels, which counts against our accuracy. An example of the mismatch between the geometry and semantics (ground truth label) is shown in the last column of Fig. 7(b), where the *sidewalk* and *road* share the same geometric plane and thus our segmentation model treats both classes as navigable space (blue color). However, according to the ground truth label, only *road* is defined as navigable.

Second, the number of degrees of freedom of the proposed method is less than that of the baselines. The degrees of freedom here can be interpreted as the dimensionality of the prediction space. Our method predicts a set of 2D boundary point coordinates, so the degrees of freedom depends on the number of points, which is usually a small number — 30 for KITTI, and 50 for Cityscapes. However, for pixel-wise methods the number of degrees of freedom is the number of pixels in the image, which is usually an enormous number — $1242 \times 375$ for KITTI and $1024 \times 512$ for Cityscapes. From the above analysis, we can see that the number of degrees of freedom of the baselines is significantly larger than that of our method, so they have higher capacity to predict highly complicated structures. However, as shown in the baseline predictions in Figs. 7(a) and 7(b), the disadvantage of many degrees of freedom is that more noise can be generated, which

| | KITTI → ORFD | | |
|---|---|---|---|
| **%** | **S**NE | **O**FF | **P**SV |
| 0 | 0.125±0.009 | 0.198±0.023 | 0.762±0.038 |
| 1 | 0.253±0.051 | 0.463±0.092 | 0.781±0.092 |
| 30 | 0.422±0.109 | 0.663±0.102 | 0.788±0.109 |
| 100 | 0.673±0.142 | 0.684±0.173 | 0.851±0.132 |

TABLE I: Comparison of domain adaptive performance for different methods.



Fig. 9: Two scenarios of indoor data, each of which is listed as one row.

might confuse planning modules during navigation.

### F. Generalization Evaluation

Generalization performance is crucial for applying deep models in real-world problems. To compare the generalization of different methods, we conduct an experiment in which all models are trained on KITTI but tested on ORFD, which is a very different dataset. We report quantitative results (measured with mIoU) in Table I. We see that the performance of baselines (SNE and OFF) are significantly degraded at all ground truth label amounts compared to the KITTI testing experiments (see Fig. 8(a)). However, our PSV-Net still maintains relatively good results.

### G. Indoor Experiments and Navigation

The experiments in the previous sections showed that our proposed PSV-Net can work well even if no ground truth labels are provided. This is a significant advantage when deploying the model to novel environments, since collecting data for new scenarios can be prohibitively expensive. We validate the proposed model on our INDOOR dataset, which has no ground truth labels. Qualitative results are illustrated in Fig. 9. The trained model on the INDOOR data can be directly combined with our visual planner to build a navigation system. Thanks to the efficient design of PSV-Net, during inference, we only use the trained $E_2$ to directly predict boundary vertices from input surface normal images. The inference speed is 36 frames per second on average, permitting real-time applications. By integrating the proposed PSV-Net with the visual planner, we succeed in building a visual navigation system on a ground robot. We only use an RGB-D camera on the robot as the perception sensor. Navigation demonstrations are provided in the supplementary video (also public video link: https://www.youtube.com/watch?v=LAliV-rUYGE).



Fig. 10: Data flow in PSV-Net. (a) Surface normal image, obtained based on the lidar data provided in KITTI dataset, hence half of values are missing. (b) Ground truth label. (c) Prior for Net-I. (d) Latent representation in Net-I. (e) Pseudo label generated by Net-I. (f) Segmentation output rendered from $D_2$. The blue region indicates freespace.

| | KITTI | | | Cityscapes | | |
|---|---|---|---|---|---|---|
| | I | II | Gain | I | II | Gain |
| **A** | 92.2 | 93.9 | +1.7 | 82.3 | 84.3 | +2.0 |
| **P** | 89.3 | 91.0 | +1.7 | 73.9 | 74.3 | +0.4 |
| **R** | 78.4 | 78.5 | +0.1 | 72.6 | 75.7 | +3.1 |
| **F** | 82.1 | 83.6 | +1.5 | 72.9 | 73.2 | +0.3 |
| **I** | 70.7 | 72.1 | +1.4 | 57.8 | 58.4 | +0.6 |

TABLE II: Quantitative comparison of Net-I and Net-II for both KITTI and Cityscapes. The columns starting from I are results for Net-I while II is for Net-II. The columns starting from Gain are advantageous gains of Net-II over Net-I.

### H. Ablation Studies

We first show the data flow in a fully self-supervised trained PSV-Net and validate the necessity of each module in our network structure (see Fig. 10) where one image from KITTI dataset is used as an example. We construct the prior using 10 randomly-selected ground truth labels in KITTI. The goal of specifying a prior is to inform the model that the navigable space is generally in the lower middle part of a first-person view image. Fig. 10(d) shows the first channel of the latent representation (two channels in total since we only perform binary segmentation) of Net-I. Fig. 10(e) is the pixel-wise label prediction results from the latent representation of Net-I. Although Fig. 10(e) is able to show the shape of the navigable space, there is still abundant noise. We use this noisy prediction as the pseudo label to train Net-II, and the prediction result (Fig. 10(f)) is accurate and clean.

To further show the advantage of combining Net-I and Net-II, we quantitatively compare the performance of the two nets for both KITTI and Cityscapes. To perform a thorough ablation study, besides the m**I**oU (see Section IV-A), we also use several other metrics, including **A**ccuracy, **P**recision, **R**ecall, **F**-Score. The definitions for those additional metrics are $\mathbf{A} = \frac{n_{tp}+n_{tn}}{n_{tp}+n_{tn}+n_{fp}+n_{fn}}, \mathbf{P} = \frac{n_{tp}}{n_{tp}+n_{fp}}, \mathbf{R} = \frac{n_{tp}}{n_{tp}+n_{fn}}$, and $\mathbf{F} = \frac{2n_{tp}^2}{2n_{tp}^2+n_{tp}(n_{fp}+n_{fn})}$ As shown in Table II, the segmentation performance is systematically improved with the refinement of Net-II. This indicates our proposed polyline-based segmentation is able to reduce the noise in pixel-wise predictions.

## V. CONCLUSION

We propose a new framework, PSV-Net, to learn the navigable space in a self-supervised fashion. The proposed framework discretizes the boundary of navigable spaces into a set of vertices. Through extensive evaluations, we have validated the effectiveness of the proposed method and its advantages over the state-of-the-art fully-supervised learning baseline methods. We also validate the effectiveness of the proposed PSV-Net with a visual planning method for efficient autonomous navigation.

## REFERENCES

[1] M. Yang, K. Yu, C. Zhang, Z. Li, and K. Yang, "Denseaspp for semantic segmentation in street scenes," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3684–3692.

[2] R. Fan, H. Wang, P. Cai, and M. Liu, "Sne-roadseg: Incorporating surface normal information into semantic segmentation for accurate freespace detection," in *European Conference on Computer Vision*. Springer, 2020, pp. 340–356.

[3] M. Cordts, M. Omran, S. Ramos, T. Scharwächter, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset," in *CVPR Workshop on the Future of Datasets in Vision*, vol. 2, 2015.

[4] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1231–1237, 2013.

[5] Z. Chen, D. Pushp, and L. Liu, "CALI: Coarse-to-Fine ALIgnments Based Unsupervised Domain Adaptation of Traversability Prediction for Deployable Autonomous Navigation," in *Proceedings of Robotics: Science and Systems*, New York City, NY, USA, June 2022.

[6] C. P. Burgess, L. Matthey, N. Watters, R. Kabra, I. Higgins, M. Botvinick, and A. Lerchner, "Monet: Unsupervised scene decomposition and representation," *arXiv preprint arXiv:1901.11390*, 2019.

[7] K. Greff, R. L. Kaufman, R. Kabra, N. Watters, C. Burgess, D. Zoran, L. Matthey, M. Botvinick, and A. Lerchner, "Multi-object representation learning with iterative variational inference," in *International Conference on Machine Learning*. PMLR, 2019, pp. 2424–2433.

[8] E. Jang, S. Gu, and B. Poole, "Categorical reparameterization with gumbel-softmax," *arXiv preprint arXiv:1611.01144*, 2016.

[9] Z. Wang, D. Acuna, H. Ling, A. Kar, and S. Fidler, "Object instance annotation with deep extreme level set evolution," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7500–7508.

[10] S. Peng, W. Jiang, H. Pi, X. Li, H. Bao, and X. Zhou, "Deep snake for real-time instance segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8533–8542.

[11] M. Kass, A. Witkin, and D. Terzopoulos, "Snakes: Active contour models," *International journal of computer vision*, vol. 1, no. 4, pp. 321–331, 1988.

[12] D. Marcos, D. Tuia, B. Kellenberger, L. Zhang, M. Bai, R. Liao, and R. Urtasun, "Learning deep structured active contours end-to-end," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8877–8885.

[13] M. Chen, T. Artières, and L. Denoyer, "Unsupervised object segmentation by redrawing," *arXiv preprint arXiv:1905.13539*, 2019.

[18] D. Acuna, H. Ling, A. Kar, and S. Fidler, "Efficient interactive annotation of segmentation datasets with polygon-rnn++," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2018, pp. 859–868.

[14] M. Tschannen, O. Bachem, and M. Lucic, "Recent advances in autoencoder-based representation learning," *arXiv preprint arXiv:1812.05069*, 2018.

[15] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.

[16] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.

[17] L. Castrejon, K. Kundu, R. Urtasun, and S. Fidler, "Annotating object instances with a polygon-rnn," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5230–5238.

[19] H. Ling, J. Gao, A. Kar, W. Chen, and S. Fidler, "Fast interactive object annotation with curve-gcn," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5257–5266.

[20] J. Liang, N. Homayounfar, W.-C. Ma, Y. Xiong, R. Hu, and R. Urtasun, "Polytransform: Deep polygon transformer for instance segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[21] D. Acuna, A. Kar, and S. Fidler, "Devil is in the edges: Learning semantic boundaries from noisy annotations," in *CVPR*, 2019.

[22] M. Zhen, J. Wang, L. Zhou, S. Li, T. Shen, J. Shang, T. Fang, and L. Quan, "Joint semantic segmentation and boundary detection using iterative pyramid contexts," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[23] D. Marin, Z. He, P. Vajda, P. Chatterjee, S. Tsai, F. Yang, and Y. Boykov, "Efficient segmentation: Learning downsampling near semantic boundaries," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.

[24] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.

[25] V. Guizilini, R. Ambrus, S. Pillai, A. Raventos, and A. Gaidon, "3d packing for self-supervised monocular depth estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 2485–2494.

[26] H. Kato, Y. Ushiku, and T. Harada, "Neural 3d mesh renderer," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3907–3916.

[27] N. Wang, Y. Zhang, Z. Li, Y. Fu, W. Liu, and Y.-G. Jiang, "Pixel2mesh: Generating 3d mesh models from single rgb images," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 52–67.

[28] T. M. Howard and A. Kelly, "Optimal rough terrain trajectory generation for wheeled mobile robots," *The International Journal of Robotics Research*, vol. 26, no. 2, pp. 141–166, 2007.

[29] T. M. Howard, C. J. Green, A. Kelly, and D. Ferguson, "State space sampling of feasible motions for high-performance mobile robot navigation in complex environments," *Journal of Field Robotics*, vol. 25, no. 6-7, pp. 325–345, 2008.

[30] C. Min, W. Jiang, D. Zhao, J. Xu, L. Xiao, Y. Nie, and B. Dai, "Orfd: A dataset and benchmark for off-road freespace detection," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 2532–2538.

[31] J. Fritsch, T. Kuehnl, and A. Geiger, "A new performance measure and evaluation benchmark for road detection algorithms," in *16th International IEEE Conference on Intelligent Transportation Systems (ITSC 2013)*. IEEE, 2013, pp. 1693–1700.

[32] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[33] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *Advances in neural information processing systems*, vol. 32, 2019.