

Semantic Segmentation of 3D LiDAR Data in Dynamic Scene Using Semi-supervised Learning

Jilin Mei, *Member, IEEE*, Biao Gao, *Member, IEEE*, Donghao Xu, *Member, IEEE*, Wen Yao, *Member, IEEE*, Xijun Zhao, *Member, IEEE*, and Huijing Zhao, *Member, IEEE*,

Abstract—This work studies the semantic segmentation of 3D LiDAR data in dynamic scenes for autonomous driving applications. A system of semantic segmentation using 3D LiDAR data, including range image segmentation, sample generation, inter-frame data association, track-level annotation and semi-supervised learning, is developed. To reduce the considerable requirement of fine annotations, a CNN-based classifier is trained by considering both supervised samples with manually labeled object classes and pairwise constraints, where a data sample is composed of a segment as the foreground and neighborhood points as the background. A special loss function is designed to account for both annotations and constraints, where the constraint data are encouraged to be assigned to the same semantic class. A dataset containing 1838 frames of LiDAR data, 39934 pairwise constraints and 57927 human annotations is developed. The performance of the method is examined extensively. Qualitative and quantitative experiments show that the combination of a few annotations and large amount of constraint data significantly enhances the effectiveness and scene adaptability, resulting in greater than 10% improvement.

Index Terms—3D LiDAR data, semantic segmentation, semi-supervised learning.

I. INTRODUCTION

Scene understanding is crucial for the safe and efficient navigation of autonomous vehicles in complex and dynamic environments, and semantic segmentation is a key technique. 3D LiDAR has been used as one of the main sensors in many prototyping systems for fully autonomous driving [1]. Semantic segmentation using 3D LiDAR data is illustrated in Fig. 1, where given a frame of input data (a), the problem is to find a meaningful label (i.e., object class in this research) for each pixel, super-pixel or region of the data (b). As 3D LiDAR is a 2.5D sensing of the surroundings, it can be represented equivalently in the form of a range image (c)-(d) in the polar coordinate system, and the problem of semantic segmentation can be solved by using either 3D points or range images as the input.

Semantic segmentation using 3D LiDAR data from outdoor scenes has been studied since the past decade [1]–[3]. The traditional process in these works [4], [5] includes the following steps: (1) preprocessing to divide a whole dataset

This work is supported by the National Key Research and Development Program of China (2017YFB1002601) and the NSFC Grants (61573027).

J. Mei, B. Gao, D. Xu and H. Zhao are with the Peking University, with the Key Laboratory of Machine Perception (MOE), and also with the School of Electronics Engineering and Computer Science

W. Yao and X. Zhao is with China North Vehicle Research Institute, Beijing, China.

Correspondence: H. Zhao, zhaohj@cis.pku.edu.cn.

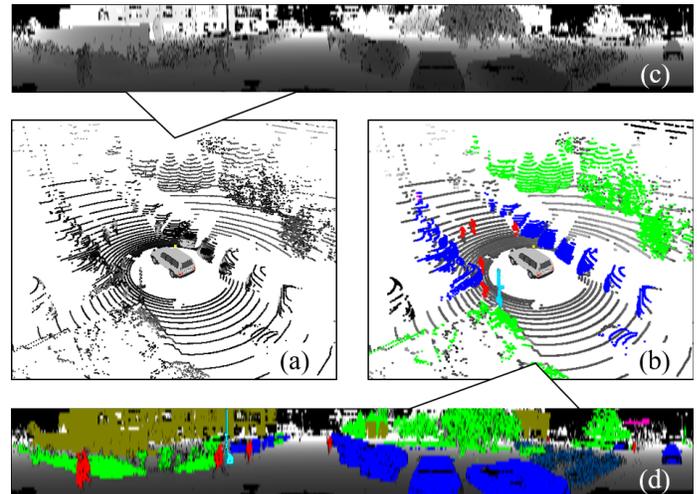


Fig. 1. The semantic segmentation for dynamic scene. (a) and (c) show the input data in two kinds of formats, i.e., the raw 3D point clouds and range frame. (b) and (d) show the semantic segmentation results.

into locally consistent small units, such as voxels, segments or clusters; (2) extracting a sequence of predefined features; (3) learning a classifier via SVM, random forest etc.; and (4) refining the results using a method such as conditional random field by considering the spatial consistency among neighboring units. The traditional methods depend on carefully designed discriminative features, and their adaptability to different scenes remains an open challenge.

The recent success of deep learning in image semantic segmentation has provided new approaches [6]. These methods remove the dependence on handcrafted features in an end-to-end manner. However, these methods also have substantial demands for finely labeled data [6]. On one hand, pixel-wise annotation is extremely time consuming; on the other hand, few 3D LiDAR datasets with annotation at the point level to support autonomous driving applications are publicly available. Therefore, it is necessary to develop a semantic segmentation method using 3D LiDAR data with only a small set of supervised data, where semi-supervised learning is adopted.

Semi-supervised learning methods, which integrate labeled and unlabeled data, have been studied extensively in the field of machine learning [7]. In [8], [9], pairwise constraints are collected from unlabeled data to describe the probability of two samples sharing the same or different labels. LiDAR sensors measure the 3D coordinates of an object directly

and can be used to associate the same object in the data of subsequent frames according to their locations after ego motion compensation. A tracking method can be used for data association for moving objects such as people, cyclists and cars. Such associated data constitute pairwise constraints, which can be inexpensive due to their abundant nature and autonomous generation.

In this paper, we propose a semantic segmentation using 3D LiDAR data from dynamic urban scenes by integrating semi-supervised learning and deep learning methods. A CNN-based classifier is trained by considering both supervised samples with manually labeled object classes and pairwise constraints, where a data sample is composed of a segment as the foreground and the neighborhood points as the background. A system of semantic segmentation using 3D LiDAR data, including range image segmentation, sample generation, inter-frame data association, track-level annotation and semi-supervised learning, is developed. A dataset containing 1838 frames of LiDAR data, 39934 pairwise constraints and 57927 human annotations is generated using 3D LiDAR data from a dynamic campus scene. Qualitative and quantitative experiments show that the combination of a small amount of annotation data and a large amount of constraint data significantly improves the effectiveness and scene adaptability of the classifier.

The remainder of this paper is organized as follows. Related work is discussed in Sect. II. In Sect. III, the proposed method is presented. In Sect. IV, the algorithm details are discussed. Then, we present the experimental results in Sect. V and draw conclusions in Sect. VI.

II. RELATED WORKS

Numerous studies on semantic segmentation have been conducted; the recent progress for image and RGB-D data is reviewed in [6]. In this section, we focus on methods specified for 3D point clouds (i.e., from LiDAR sensors) in dynamic outdoor scenes. These works can be broadly divided into three classes: feature-based methods, deep learning methods, semi-supervised learning methods.

A. Feature-based Methods

Feature-based methods belong to traditional machine learning, and the general process of these methods consists of feature selection, classifier design and graphical model description.

A straightforward technique is to convert semantic segmentation into a point-wise classification that includes extracting the features on each unit, concatenating the features as a vector and determining the label via a well-trained classifier. [10] presents a common pipeline from feature selection to classifier training. Due to the irregular arrangement of point clouds, the authors test 5 definitions of neighborhood to achieve the best representation. Similar research is conducted in [13], which demonstrates the ability to address varying densities of data. A single point cloud usually contains millions of points, so evaluating the label for each point is typically computationally expensive (on the order of minutes according to [13]).

[14] proposes an efficient approach where speed and accuracy are satisfied simultaneously; furthermore, the average classification time can be reduced to less than 1 s. [5] represents the raw point cloud as a 2D range image and proposes a framework for simultaneous segmentation and classification of the range image that considers both the 2D range image and 3D raw data. Straightforward approaches assume that each data unit is independent and ignore the spatial and contextual relations between units. Consequently, they can produce good results based on distinctive features. However, when the features are not discriminative, the point-wise classification will be noisy and locally inconsistent [4].

The neighbor elements are taken into account to make the segmentation results spatially smooth. For this purpose, graphical models such as Markov random Field (MRF) and conditional random field (CRF) are usually exploited to encode the spatial relationships. In [11], the node potentials and edge potentials are both formulated with a parametric linear model, and the functional max-margin learning is used to find the optimal weights. [16] proposes a simplified Markov network to infer the contextual relations between points. Instead of learning all the weights for the node and edge potentials in graphical models, the node potentials are calculated from a point-wise classifier, and the edge potentials are determined by the physical distance between points.

The performance of the above methods largely depends on handcrafted features. These methods are effective in fixed or regular scenarios, but for dynamic scenes, the features are empirically designed and the performance decreases.

B. Deep Learning Methods

Deep learning, especially the convolution neural network (CNN) without handcrafted features, has shown effectual performance on 2D image segmentation [6]. However, the semantic segmentation of 3D point clouds (i.e., from LiDAR sensors) is still an open research problem [17] due to the irregular, not grid-aligned properties. Therefore, recent studies project the point clouds into 2D views, and some of them attempt to directly ways, for example, volumetric/voxel representations.

Inspired by the success of CNN in image segmentation, the state-of-the-art image-based algorithms can be used directly after rendering 2D views from the 3D raw data. [20] projects point clouds into virtual 2D RGB images via Katz projection. Then, a pretrained CNN is used to semantically classify the images. However, this projection removes all the points that are not visible; for example, if a car is projected, all the points behind it are removed. [22] unwraps 360° 3D LiDAR data onto a spherical 2D plane without point loss. Spherical projection is also applied in SqueezeSeg [25], where the CNN directly outputs the point-wise label of the transformed LiDAR data and a CRF is applied to refine the outputs. [26] uses cylindrical projection to create the depth and reflectivity images. In [27], the point clouds are encoded by top-view images and a simple fully convolutional neural network (FCN) is used. This method can be used in real time because only elevation and density features are extracted. In [28], the input point cloud is projected into multiple views, such as color, depth and surface normal images.

TABLE I
APPROACHES FOR 3D POINT CLOUDS SEMANTIC SEGMENTATION

Research	Learning Method	Input	Classifier	Dataset	Scene
Munoz, 2009, [4]	sup.	L	MRF	-	rural
Zhao, 2010, [5]	sup.	L	SVM	-	campus
Weinmann, 2014 [10]	sup.	L	RF	VMR-Oakland [11],Pairs-rue-Madame [12]	urban
Munoz, 2009, [11]	sup.	L/V	MRF	-	-
Hackel, 2016 [13]	sup.	L	RF	Pairs-rue-Madame	urban
Hu, 2013, [14]	sup.	L	KLR	VMR-Oakland,Freiburg [15]	urban
Lu, 2012, [16]	sup.	L	CRF,SVM	VMR-Oakland	urban
Engelmann, 2017, [17]	sup.	L/V	DL	S3DIS [18],vKITTI [19],KITTI	indoor,urban,urban
Tosteberg, 2017, [20]	sup.	L	DL	Semantic3D,VPS [21]	urban,indoor
Dewan, 2017, [22]	sup.	L	DL	KITTI [23]	urban
Hackel, 2017, [24]	sup.	L	DL	Semantic3D [24]	urban/rural
Wu, 2017, [25]	sup.	L	DL	KITTI	urban
Piewak, 2018, [26]	sup.	L/V	DL	-	urban/rural/highway
Caltagirone, 2017, [27]	sup.	L	DL	KITTI	urban
Lawin, 2017, [28]	sup.	L	DL	Semantic3D	urban/rural
Tchapmi, 2017, [29]	sup.	L/V	DL	NYU V2 [30],Semantic3D,S3DIS,KITTI,	indoor,urban/rural,indoor,urban
Riegler, 2017, [31]	sup.	L	DL	ModelNet10 [32]	CAD model
Qi, 2017, [33]	sup.	L/V	DL	ShapeNet [34],Stanford3D [35]	CAD model,indoor
Landrieu, 2017, [36]	sup.	L/V	DL	Semantic3D,S3DIS	urban,rural/indoor
Bearman, 2016, [37]	semi-sup.	V	DL	PASCAL VOC [38]	-
Yan, 2006, [39]	semi-sup.	V	KLR,SVM	-	nursing home
Bauml, 2013 [40]	semi-sup.	V	KLR,SVM	-	TV series
Hong, 2015, [41]	semi-sup.	V	DL	PASCAL VOC	-
Papandreou, 2015, [42]	semi-sup.	V	DL	PASCAL VOC	-
Lin, 2016, [43]	semi-sup.	V	DL	PASCAL VOC	-
Cour, 2009, [44]	weakly-sup.	V	linear classifier	-	TV series
Pathak, 2015, [45]	weakly-sup.	V	DL	PASCAL VOC	-
Xu, 2015, [46]	weakly-sup.	V	linear classifier	Siftflow [47]	-
Dai, 2015, [48]	weakly-sup.	V	DL	PASCAL VOC	-

sup. : supervised; L : LiDAR Data; V : Camera Data; RF : Random Forest; KLR : Kernel Linear Regression; DL : Deep Learning.

Another type of method models the raw data in direct ways. [29] proposes SEGCloud, where the raw 3D point cloud is preprocessed into a voxelized point cloud with a fixed grid size. Although [29] is simple and effective, how to set the voxel size is a problem in large-scale scenes. Thus, the scene is voxelized at five different resolutions in [24], and each of the five scales is handled separately by the CNN. Rather than using a fixed grid size, [31] proposes OctNet, where the hybrid grid-octree data structure is applied to represent the raw 3D data, and each leaf of the octree stores a pooled feature representation. PointNet [33] is a unified architecture that directly takes raw point clouds as input and outputs the label of each point. The scene is divided into blocks. Then, the points in each block are passed through a series of multilayer perceptrons (MLPs) to extract the local and global features. Based on [33], [17] extends the method to incorporate a larger-scale spatial context, and improved results are reported in both indoor and outdoor scenarios. [36] proposes a more elegant architecture to capture contextual relations. The first step is to

partition the raw point cloud into geometrically simple shapes, called super-points. The super-points are then embedded by PointNet [33].

Semantic segmentation with deep learning is usually implemented in a supervised manner, which requires detailed annotations. However, obtaining point-wise annotations for 3D point clouds is labor intensive and time consuming. Furthermore, few public datasets support this level of annotation.

C. Semi-supervised Learning Methods

Considering the large demand for detailed annotations, many researchers study semi-supervised learning methods, which integrate fewer labeled data and more unlabeled data, and weakly supervised learning methods, which use multiple ambiguous labels. Our research belongs to the semi-supervised category, please refer to TABLE I for weakly supervised methods.

The early work [8] on semi-supervised learning specifies the prior knowledge of unlabeled data via pairwise constraints. A

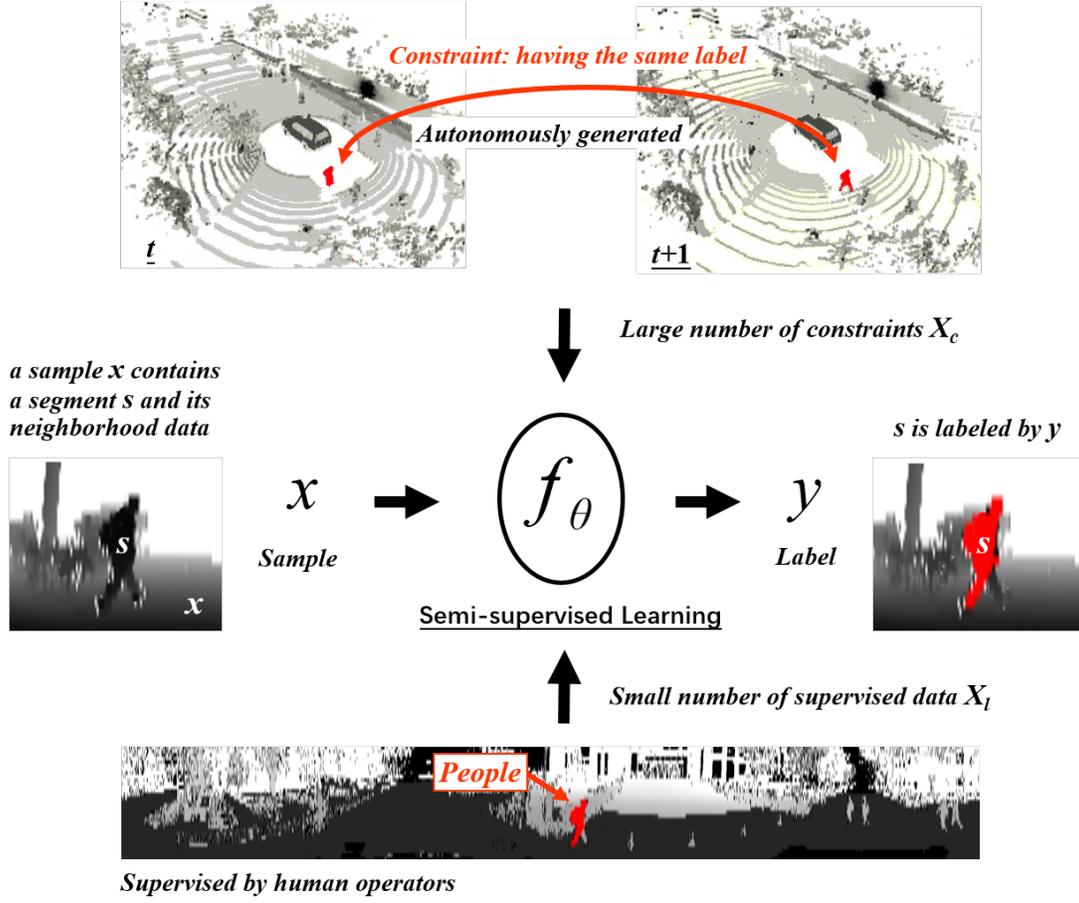


Fig. 2. The framework of semi-supervised learning for 3D point clouds semantic segmentation.

pairwise must-link constraint means two objects must have the same label, although the label is unknown, whereas two objects associated via must-not-link must have different labels. Both labeled and constraint data are used for model fitting, and the authors model the constraint information by the maximum entropy principle. A similar idea is presented in [9]. The pairwise constraints are incorporated in the clustering of a Gaussian mixture model (GMM). [8], [9] present the foundation of semi-supervised learning with pairwise constraints.

[39] extends the method for video object classification, where temporal relations between frames, multi-modalities, such as faces and voices, and human feedback (manual annotation) are considered and formulated in a unified framework. [40] applies constraints for person identification in multimedia data and achieves state-of-the-art performance on two diverse TV series. Recently, semi-supervised learning has also been integrated with deep neural networks (DNN). [41] decouples semantic segmentation into classification and segmentation using a large number of image-level object labels and a small number of pixel-wise annotations. The classification network specifies the class-specific activation maps, which are transferred into the segmentation network with bridge layers. Then, the segmentation network requires only a few pixel-wise annotations to train the model, e.g., 5 or 10 strong annotations per class. [42] designs a expectation-maximization

(EM) training method for semantic image segmentation by combining bounding boxes, image-level labels and a few strongly labeled images.

Semi-supervised learning has been successfully applied in image segmentation [41] and video analysis [39]. However, to the best of the author's knowledge, few studies discuss semi-supervised learning in the context of the 3D point clouds in semantic segmentation. In this paper, we attempt to combine semi-supervised learning and neural network to solve the problem of how to perform 3D point cloud semantic segmentation with insufficient point-wise annotations.

III. METHODOLOGY

A. Problem Definition

Let s denotes a small segment of 3D LiDAR data extracted by examining the consistency of 3D points with their neighborhood in the range image frame using, e.g., a region growing method. Without loss of generality, we assume that the 3D points of s are measurements of a single object. However, s commonly represents only a part of the object, e.g., the upper body of a pedestrian, due to the nature of LiDAR measurements. Hence, for each s , a data sample x is generated centered at s , containing s as the foreground and its neighborhood data as the background, as shown in Fig. 2.

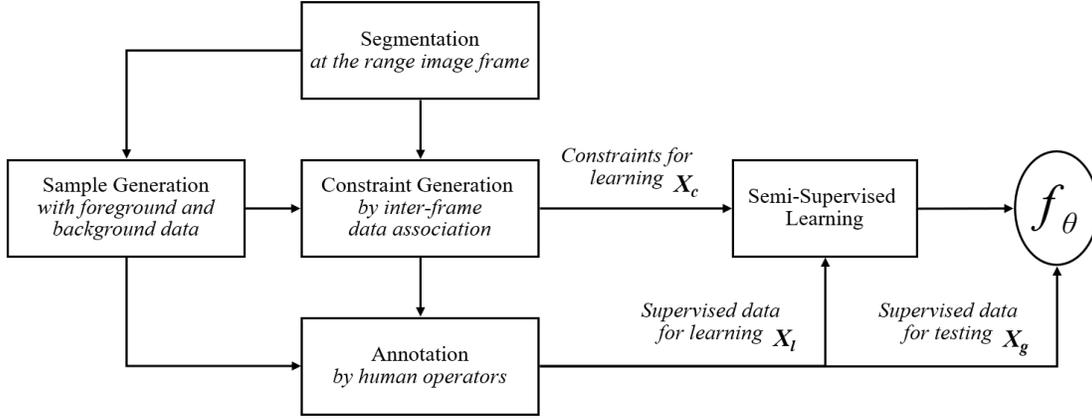


Fig. 3. The flowchart of implementation.

The problem in this work is formulated as learning a multi-class classifier f_θ that maps x to a label $y \in \{1, \dots, K\}$ and subsequently associates y with the 3D points of s .

$$f_\theta : x \rightarrow y \in \{1, \dots, K\} \quad (1)$$

Given a set of supervised data samples $X_l = \{x_i, y_i\}_{i=1}^M$, where $\{y_i\}$ are one-hot vectors annotated manually by human operators for each $\{x_i\}$, a common way of learning a classifier f_θ is to find the best θ^* that minimizes a loss function L , as below.

$$\theta^* = \arg \max_{\theta} L(X_l; \theta) \quad (2)$$

However, the problem of generating a large amount of supervised data is not trivial. This research learns f_θ with a small set of costly supervised data X_l and an additional large set of autonomously generated constraints $X_c = \{\langle x_i, x_j \rangle\}_{n=1}^N$,

$$\theta^* = \arg \max_{\theta} L(X_l, X_c; \theta) \quad (3)$$

where constraint $\langle x_i, x_j \rangle$ denotes that x_i and x_j have the same label, i.e., $y_i = y_j$, which is generated in this research autonomously by associating the data segments along sequential frames according to their locations after ego-motion compensation.

This semi-supervised loss function L is based on a combination of loss on supervised data X_l and loss on constraints X_c .

$$L(X_l, X_c; \theta) = L_l(X_l; \theta) + L_c(X_c; \theta) \quad (4)$$

B. Supervised Loss

For supervised data X_l , we follow the widely used definition of cross entropy, e.g., [33], and define loss L_l as below:

$$L_l(X_l; \theta) = -\frac{1}{M} \sum_{i=1}^M \sum_{k=1}^K \mathbf{1}[y_i^k = 1] \ln(P_\theta^k(x_i)) \quad (5)$$

where $\mathbf{1}[*]$ is an indicator function, and $P_\theta^k(x_i)$ is the probability that x_i is assigned a label k by a classifier with the set of parameters θ .

C. Constraint Loss

For each constraint $\langle x_i, x_j \rangle$, let y_i and y_j denote the labels x_i and x_j , respectively, by using a learned classifier f_θ . A penalty is applied if $y_i \neq y_j$. Subsequently, the loss is estimated as below for each constraint based on the probability that the constrained data samples are assigned different labels.

$$\begin{aligned} P(y_i \neq y_j) &= \sum_{k=1}^K \sum_{\substack{l=1 \\ l \neq k}}^K P_\theta^k(x_i) P_\theta^l(x_j) \\ &= 1 - \sum_{k=1}^K P_\theta^k(x_i) P_\theta^k(x_j) \end{aligned} \quad (6)$$

Hence, we define the loss on constraints L_c as below:

$$\begin{aligned} L_c(X_c; \theta) &= \frac{\gamma}{N} \sum_{n=1}^N P(y_i \neq y_j) \\ &= \frac{\gamma}{N} \sum_{n=1}^N \left(1 - \sum_{k=1}^K P_\theta^k(x_i) P_\theta^k(x_j)\right), \gamma \in [0, 1] \end{aligned} \quad (7)$$

where γ is a weighting factor. Clearly, L_c describes the unsupervised learning procedure.

D. Semi-supervised Loss

Combining the losses from supervised data X_l and constraints X_c , the semi-supervised loss is defined as below.

$$\begin{aligned} L(X_l, X_c; \theta) &= L_l(X_l; \theta) + L_c(X_c; \theta) \\ &= -\frac{1}{M} \sum_{i=1}^M \sum_{k=1}^K \mathbf{1}[y_i^k = 1] \ln(P_\theta^k(x_i)) + \\ &\quad \frac{\gamma}{N} \sum_{n=1}^N \left(1 - \sum_{k=1}^K P_\theta^k(x_i) P_\theta^k(x_j)\right), \gamma \in [0, 1]. \end{aligned} \quad (8)$$

IV. ALGORITHM DETAILS

A. Process Flow

Fig. 3 describes the major modules of the workflow. Sample generation and semi-supervised learning are detailed in the next subsections. Here, we describe the remaining modules. Although traditional methods are used in these modules, their

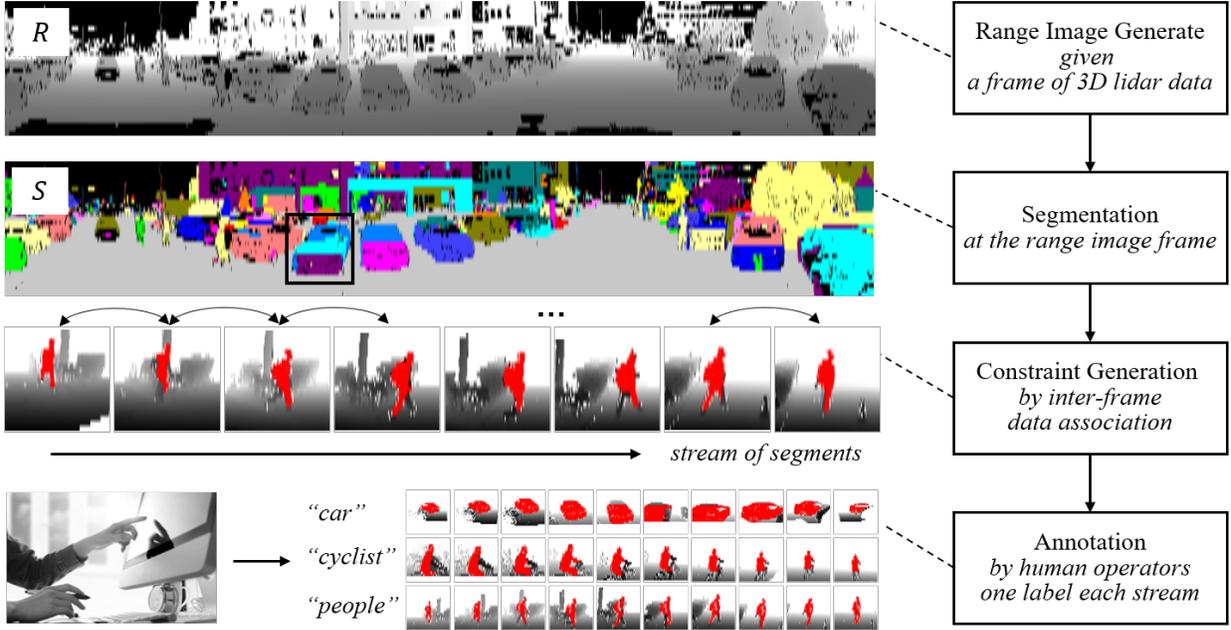


Fig. 4. The details for segmentation, constraint generation and annotation. The first row is range image; the second row is the segmentation result and the black rectangle shows a car is separated into three parts; the third row shows two adjacent sample consist of one constraint; the fourth row is human annotation.

integration and the design of the data pipelines are important in constructing a complete system.

As illustrated in Fig. 4, segmentation is conducted at the frame level of a range image, which is a representation of 3D LiDAR data in the polar coordinate system. The columns and rows of a range image correspond to tessellated horizontal and vertical angles, and each pixel value is the range distance of the laser point in that direction. A Velodyne HDL-32E is used in this research. Hence, a LiDAR frame contains 32 scan lines at different vertical angles, and each scan line has 2160 laser points at different horizontal angles. These laser points are projected onto a range frame and reshaped to size (144,1080).

Region growing is conducted to extract segments from unlabeled pixels as the seeds by examining 4-connectivity, where the thresholds on the vertical and horizontal range differences of two connected pixels are assigned empirically with a set of test data. As detailed in the next subsection, a segment is used as the foreground data in sample generation, and segmentation is treated only as a preprocessing step. Many methods can be exploited as long as the following condition is met: the 3D points of segment s are measurements of a single object.

Constraints are generated by inter-frame data association. For any segment s_t in frame t , the 3D points can be back-projected to the LiDAR sensor’s coordinate system in the previous frame $t - 1$ based on the vehicle’s motion data and calibration parameters. A segment s_{t-1} is associated with s_t if it matches with the 3D points. Let x_i and x_j denote the samples of s_{t-1} and s_t , respectively; a pairwise-constraint $\langle x_i, x_j \rangle$ is subsequently generated.

As illustrated in Fig. 4, data association is conducted for the entire stream, and sequences of segments are extracted as the result. An operator examines each sequence. If the sequence is

tracked correctly, a label is assigned to all the segments (and the samples) of the sequence, and constraints are generated for all subsequent samples. Otherwise, a sequence is manually truncated to remove the erroneous tracking part or even discarded.

The dataset with supervised labels is divided into two groups, X_l and X_g , for training and testing, respectively. The set of pairwise constraints X_c is used for training only.

B. Sample Generation

A straightforward method is to use a segment s directly as the input of a classifier, i.e., $s = x$. However, the performance of such a classifier can be degraded if s represents only a part of the target object. This situation often occurs in the segmentation results of 3D LiDAR data. Due to diffusive reflection or the weak reflectance of LiDAR measures on reflective or dark objects, many failures that yield discontinuities in the data of a single object exist in 3D LiDAR measurements. As indicated by the black rectangle in the second row of Fig. 4, three segment pieces are extracted from the data of a single car, and it is difficult to recognize the car given the data of only one piece. However, by placing each piece of the segment into the background of its surrounding data, the car is easily recognized.

Inspired by the above idea, given a segment s , this work generates a sample x containing s as the foreground and its neighborhood data as the background. As illustrated in Fig. 5(b), a cuboid centered at s with a size of 2.4 m x 5 m x 5 m (height x width x length) is drawn. The LiDAR points inside the cuboid are extracted, and their pixels on the range image are projected onto a canvas with a size of 256x256 to obtain sample x , where each pixel is composed of three channels: 1) Height: distance to the ground surface mapped to

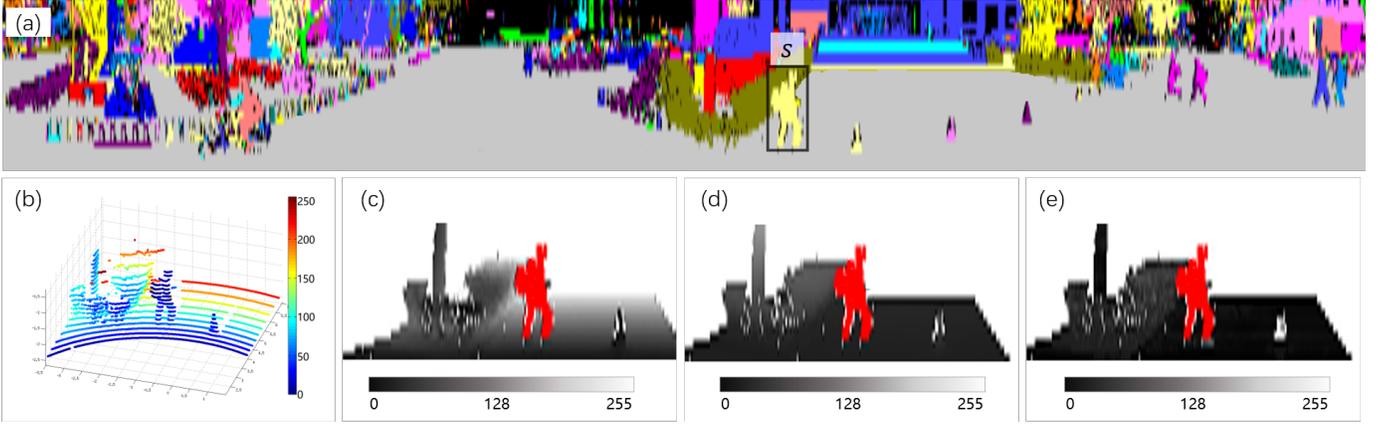


Fig. 5. The procedure of sample generation. (a) the yellow segment s in the black rectangle is chosen as candidate region. (b) the raw points inside a cuboid centered at s are cropped, where the cuboid size is $2.4\text{m} \times 5\text{m} \times 5\text{m}$ (height,width,length) and the points are colored by range value; then these points are projected on range image to make one sample that consists of three channels. (c) the range channel of the sample, and we mark s with red for better visualization. (d) the height channel of the sample. (e) the intensity channel of the sample.

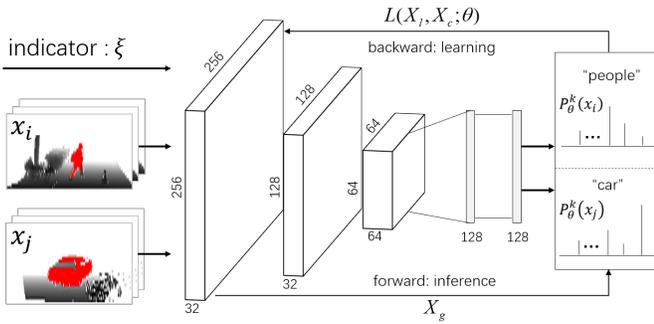


Fig. 6. The classifier for offline semi-supervised learning.

[0,255]. Distances in [0 m,6 m] are mapped linearly to [0,255]. Distance less than 0 m or greater than 6 m are mapped to 0 and 255 respectively; 2) Range: distance to the sensor, normalized to [0,255]. 3) Intensity: reflectance of the LiDAR point in [0,255].

As a small or distant segment may provide insufficient information for a reliable classification, the following criteria are applied.

$$\frac{s_n}{s_d} > \rho \cap s_n > \sigma \quad (9)$$

where s_n is the number of LiDAR points of segment s , s_d is the distance from the LiDAR sensor to the center point of s , and ρ and σ are empirically assigned thresholds. Segment s is valid for sample generation if it has more points than σ and the n-d ratio is larger than ρ . In this research, $\sigma = 8.0$ and $\rho = 30$.

C. Semi-supervised Learning

A CNN is used as the classifier f_θ , as shown in Fig. 6, with specifically designed input and loss function. We rewrite the loss function of Equation (8) as below.

$$L(X_l, X_c; \theta) = \begin{cases} L_c(X_c; \theta), & M = 0, N \neq 0 \\ L_l(X_l; \theta), & M \neq 0, N = 0 \\ L_l(X_l; \theta) + L_c(X_c; \theta), & M \neq 0, N \neq 0. \end{cases} \quad (10)$$

If no supervised samples exists, i.e., $M = 0, N \neq 0$, the loss function degenerates to $L_c(X_c; \theta)$, representing unsupervised learning. If there are no constraints, i.e., $M \neq 0, N = 0$, the loss function becomes $L_l(X_l; \theta)$, representing supervised learning. Finally, if both supervised samples and constraints exist, i.e., $M \neq 0, N \neq 0$, the loss function is in its full form, i.e., $L_l(X_l; \theta) + L_c(X_c; \theta)$, representing semi-supervised learning. The classifier is designed to adapt to all the above cases.

Algorithm 1 the training of the CNN classifier

Input: X_c, X_l

Output: the classifier parameter θ

- 1: Initialize Φ_l, Φ_c with empty.
 - 2: Make input pairs:
 - 3: $\Phi_c \leftarrow \langle x_i, x_j; \xi = 1 \rangle, \forall \langle x_i, x_j \rangle \in X_c$
 - 4: $\Phi_l \leftarrow \langle x_i, y_i, x_j, y_j; \xi = 0 \rangle, \forall \langle x_i, y_i, x_j, y_j \rangle \in X_l$
 - 5: **for each step in training do**
 - 6: $\Phi_c^n \leftarrow$ take n items from $\Phi_c, n > 0$
 - 7: $\Phi_l^m \leftarrow$ take m items from $\Phi_l, m > 0$
 - 8: $\Phi = \Phi_l^m \cup \Phi_c^n$
 - 9: **for each item in Φ do**
 - 10: **if $\xi = 0$ then**
 - 11: $L(X_l, X_c; \theta) = L_l(X_l; \theta)$
 - 12: $= L_l(x_i, y_i, x_j, y_j; \theta)$
 - 13: **else**
 - 14: $L(X_l, X_c; \theta) = L_c(X_c; \theta) = L_c(x_i, x_j; \theta)$
 - 15: Do backward learning.
 - 16: return θ
-

To allow both supervised samples and pairwise constraints in model training, the CNN is designed to take two samples x_i and x_j as input and output their labels y_i and y_j simultaneously. An indicator ξ is used to specify whether the two samples are a constrained pair ($\xi = 1$) or individuals ($\xi = 0$). Hence, the loss functions are converted to the following.

TABLE II
THE DATA SET.

	Frame	Dist.(m)		People	Car	Cyclist	Trunk	Bush	Building	Unknown	Total
A	0~414	0~184	cons.	682	1897	196	975	3268	2450	0	9468
			anno.	735	2079	228	1048	4330	2825	2423	13668
B	414~829	184~350	cons.	681	1896	195	977	3268	2450	0	9467
			anno.	736	2080	227	1048	4330	2825	2423	13669
C	829~1373	350~630	cons.	909	2529	271	1301	4357	3267	0	14624
			anno.	980	2772	304	1398	5773	3767	3230	18224
D	1373~1838	630~890	cons.	925	3542	142	62	2440	1254	0	8365
			anno.	962	4157	169	89	3773	1549	1667	12366
Total	1838	890	cons.	3197	9864	804	3315	13333	9421	0	39934
			anno.	3413	11088	928	3583	18206	10966	9743	57927

¹ The A-D corresponding the route in Fig.7.

² Dist.: the travel distance. cons.: constraint. anno.: annotation.

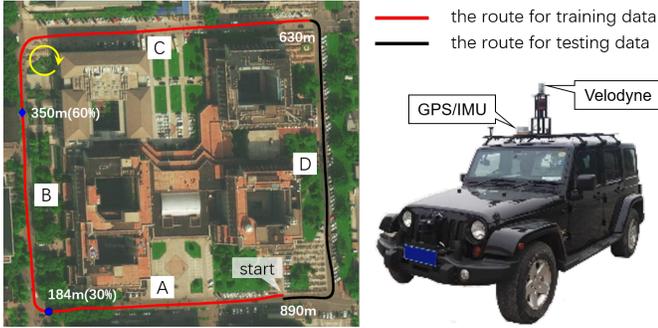


Fig. 7. The routes of data collection and the platform configuration. The 30%/60% in the left means 30%/60% travel of the route for training data.

$$\begin{aligned}
 L(X_l, X_c, \xi = 1; \theta) &= L_l(X_l; \theta) = L_l(x_i, y_i, x_j, y_j; \theta) \\
 &= -\frac{1}{2} \sum_{t=i}^j \sum_{k=1}^K \mathbf{1}[y_t^k = 1] \ln(P_\theta^k(x_t)); \\
 L(X_l, X_c, \xi = 0; \theta) &= L_c(X_c; \theta) = L_c(x_i, x_j; \theta) \\
 &= \frac{1}{2} \left(1 - \sum_{k=1}^K P_\theta^k(x_i) P_\theta^k(x_j) \right).
 \end{aligned} \tag{11}$$

In training, the supervised samples X_l and constraints X_c are randomly fed into the CNN, and the model parameters are adjusted in the traditional back-propagation manner. Pseudo code of the training process is given in Algorithm 1. As the focus of this research is to learn a classifier using small number of expensive supervised samples and a large number of inexpensive constraints, the unbalanced sample problem should be considered. In this study, we set the number of constraints to 5 times the number of supervised samples, which is described as $n/m = 5$ in lines 6-7 of Algorithm 1.

In online classification, we have $\xi \equiv 0$, and in the case of only one sample, we have $x_i = x_j$.

V. EXPERIMENTAL RESULTS

A. Data Set

The performance of the proposed method is evaluated on a dynamic campus dataset collected by an instrumented vehicle with a GPS/IMU suite and a Velodyne-HDL32, as shown in Fig. 7. The total route is approximately 890 meters. All sensor data are collected, and each data frame is associated with a time log for synchronization. The GPS/IMU data are logged at 100 Hz. The LiDAR data are recorded at 10 Hz and include 1373 frames of training data (red line in Fig. 7) and 465 frames of testing data (black line in Fig. 7). One frame can produce multiple samples, for example, we obtain 6931 car samples from the 1373 frames of training data in TABLE II. In total, we obtain 1838 frames of LiDAR data, 39934 constraints and 57927 manual annotations.

The details of the dataset are listed in TABLE II. Six labels are used, i.e., person, car, cyclist, trunk, bush and building. These categories are important for driving applications on a campus; other labels, such as pole and cone, are marked as unknown. We do not generate constraints for the unknown label.

B. Result - Classifier Training

1) *Experimental settings*: There are five experimental settings for the training data, as shown in Table III. The training data contain 3 parts, i.e., A, B, C in TABLE II. A total of 70% of the data are randomly selected for training, and the remaining 30% are selected for validation. Except baseline_max, the settings use only a small amount of annotations. Baseline_min uses 350 annotations without the constraint, while baseline_max uses all the annotations. Thus, the baseline method works in a supervised manner. In general, baseline_min sets the low performance bound and baseline_max sets the high performance bound. Constraint30 uses the same annotations as the baseline_min but with additional constraints; tail 30 means all constraints are used during 30% of the travel, i.e., route A from the start to the 184 m point (30%) in Fig. 7.

TABLE III
THE EXPERIMENTAL SETTINGS ON TRAINING DATA.

Settings		People	Car	Cyclist	trunk	Bush	Building	Unknown
baseline_min	constraint	0	0	0	0	0	0	0
	annotation	350	350	350	350	350	350	5650
baseline_max	constraint	0	0	0	0	0	0	0
	annotation	1715	4851	531	2445	10103	6591	5650
constriant30	constraint	682	1897	196	976	3268	2450	0
	annotation	350	350	350	350	350	350	5650
constriant60	constraint	1363	3793	391	1952	6536	4900	0
	annotation	350	350	350	350	350	350	5650
constriant100	constraint	2272	6322	652	3253	10893	8167	0
	annotation	350	350	350	350	350	350	5650

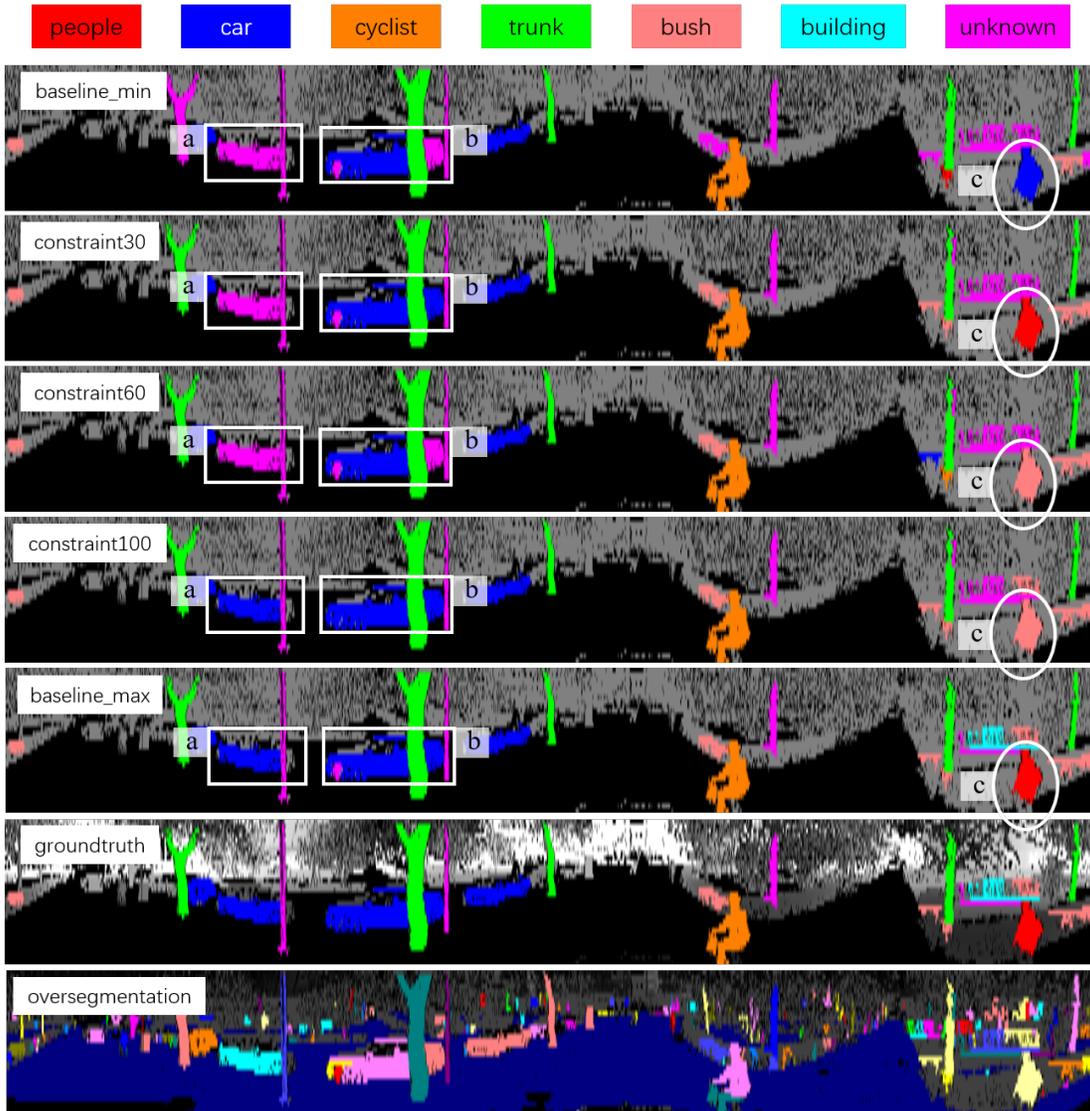


Fig. 8. The qualitative results on training data. The rectangles a and b show that the constraint100 has better performance on car, and the ellipse c shows that the constraint100 makes an error classification when the people near the bush.

TABLE IV
F MEASURE ON TRAINING DATA.

Learning Method	Classifier	People	Car	Cyclist	Trunk	Bush	Building	Unknown
sup.	baseline_min	77.2	68.5	66.1	81.0	56.2	78.7	33.9
semi-sup.	constraint30	77.2	83.0	69.6	86.6	61.9	81.8	36.8
semi-sup.	constraint60	81.5	85.8	71.3	90.4	80.7	86.6	48.3
semi-sup.	constraint100	87.0	90.4	77.8	90.6	81.2	87.7	52.2
sup.	baseline_max	93.3	96.4	83.8	96.3	92.3	94.3	80.0

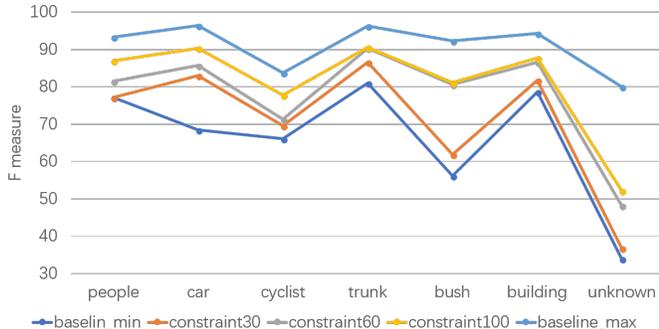


Fig. 9. The quantitative comparison on training data.

constraint60 and constraint100 have similar configurations to constraint30, except for the amount of constraint data.

According to the settings, the five classifiers are learned separately offline. All the classifiers have the same network architecture, as detailed in Set. IV.C. The difference lies in the loss function, where baseline_min and baseline_max use only $L_l(X_l; \theta)$ in Equation (8), and the other settings use both $L_l(X_l; \theta)$ and $L_c(X_c; \theta)$. For the offline learning, the classifier's parameters are saved at fixed training steps; then, the parameters that make the loss less than $1e-4$ and achieve the best performance on the validation set is selected. For online inference, all classifiers work in the same way.

2) *Qualitative results*: As illustrated in Fig. 8, as the number of constraints increases, the car in the rectangle is successfully classified by constraint100, even if occlusions occur. Our method still produces errors, for example, when the person walks near the bush, the constraint100 classifier wrongly annotates the person as a bush. The main reason for this error is that a single sample contains both the foreground and background; if the background occupies more information than the foreground, the classifier is likely to assign the background label to the sample.

3) *Quantitative results*: The F-measure is adopted for quantitative evaluations and is defined as:

$$F - Measure = \frac{2 * recall * precision}{recall + precision} \cdot 100\%. \quad (12)$$

We use the five classifiers to annotate the training set and the validation set, and the quantitative results are shown in TABLE IV.

The baseline_min and baseline_max results show that a large number of annotations is important for supervised

learning. The baseline_min and constraint100 results indicate that semi-supervised learning is effective: the F-Score of constraint100 increases by 15% on average. Furthermore, as the number of constraints increases, the performance of the classifier is enhanced. A more intuitive comparison is illustrated in Fig. 9. Although constraint100 is not as good as baseline_max, in which all annotations are used, it shows promising results, indicating that adding constraints improves the performance of the classifier and reduces the need for annotations. In conclusion, semi-supervised learning, where only a few annotations are used, is effective for 3D point cloud semantic segmentation.

C. Result - Classifier Testing

In a fixed scene or dataset, a classifier based on supervised learning can easily achieve high performance due to overfitting. When the classifier is applied to a new scene, additional annotations are necessary to prevent performance degradation. However, large quantities of fine annotations in each new scene are difficult to obtain for driving applications. Thus, how to enhance the adaptability with a few or no new annotations is crucial for practical applications. We detail three experiments in the following to demonstrate the adaptability of the proposed semi-supervised method. The F-measure in Equation (12) is used for the quantitative comparison.

1) *The pretrained result*: The pretrained classifiers based on training data are directly applied to the testing data, and the results are shown in TABLE V. The baseline_min and constraint100 results show that adding constraints improves the adaptability to the new scene: the F-Score of constraint100 increases by 9% on average. The constraint100 and baseline_max results show that constraint100 has higher scores for the cyclist and trunk categories, despite having lower performance in all categories on the training data (TABLE IV).

2) *The unsupervised result*: For the new scene, an attempt to use only constraints and no new annotations to improve the pretrained classifier is interesting. Thus, the loss function in Equation. (8) is rewritten as:

$$L(X_l, X_c; \theta) = L_c(X_c; \theta) = \frac{1}{N} \sum_{i=1}^N \left(1 - \sum_{k=1}^K P_{\theta}^k(x_i) P_{\theta}^k(x_j) \right). \quad (13)$$

Here, the pretrained constraint100 is treated as the initial classifier, and only constraints are used to retrain the model.

TABLE V
THE PRE-TRAINED CLASSIFIERS ON TESTING DATA.

Learning Method	Classifier	People	Car	Cyclist	Trunk	Bush	Building	Unknown
sup.	baseline_min	57.6	64.0	26.3	37.5	33.2	49.7	33.0
semi-sup.	constraint100	69.4	81.6	39.8	42.1	37.3	55.7	39.0
sup.	baseline_max	73.7	88.8	30.0	30.8	71.0	65.4	53.0

TABLE VI
THE SETTING OF FINE TUNING

	People	Car	Cyclist	Trunk	Bush	Building	Unknown
anchor sample	20	20	20	20	20	20	100
constraint	925	3542	142	62	2440	1254	0

TABLE VII
F MEASURE ON TESTING DATA

Learning Method	Classifier	People	Car	Cyclist	Trunk	Bush	Building	Unknown
sup.	baseline_min	57.6	64.0	26.3	37.5	33.2	49.7	33.0
semi-sup.	constraint100	69.4	81.6	39.8	42.1	37.3	55.7	39.0
semi-sup.	constraint100+fine tuning	77.1	90.7	60.2	55.7	58.0	62.3	54.3
sup.	baseline_max	73.7	88.8	30.0	30.8	71.0	65.4	53.0

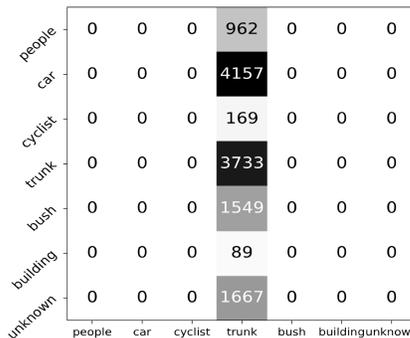


Fig. 10. The confusion matrix of unsupervised learning on testing data.

The results of this unsupervised learning procedure are shown in Fig. 10. Regardless of the samples, the classifier erroneously assigns them as a trunk. We can explain this situation with the loss function in Equation (13): the constraint loss only penalizes the classifier when it gives x_i and x_j different labels, so unsupervised learning fails.

3) *The fine tuning result:* After finding that unsupervised learning does not work for our method, we attempt to add a few new annotations. Specifically, the new annotations are generated in an interactive way. The pretrained constraint100 is used as the initial classifier to produce classification results on the testing data. Although the pretrained results show low F-Scores in TABLE V, a few new annotations can be selected by human confirmation. In this way, we obtain 20 annotations for each category and 100 for the unknown label, as shown in TABLE VI, where the new annotation is renamed the anchor sample.

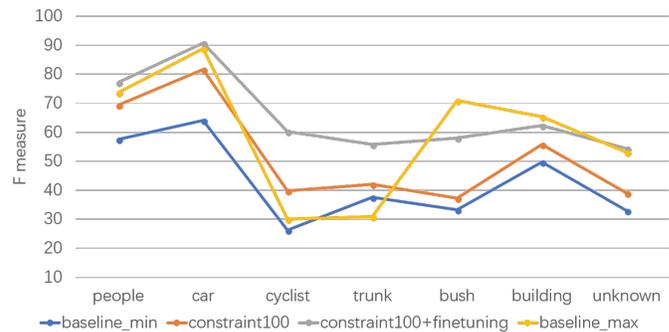


Fig. 11. The quantitative comparison on testing data.

The pretrained constraint100 model is fine-tuned by combining the anchor sample and constraint. The final quantitative results are shown in TABLE VII and Fig. 11. Comparing constraint100 and the fine-tuned version, the F-Score of the latter increases by 13% on average, which shows that fine-tuning is an effective way to improve adaptability, even in a semi-supervised manner. Comparing the baseline_max and the fine-tuned version, the latter has higher scores except on bush and building. The qualitative results are shown in Fig. 12. In conclusion, fine-tuning with the anchor sample increases the adaptability of the classifier to a new scene.

VI. CONCLUSION AND FUTURE WORK

A semantic segmentation method for 3D point clouds (i.e., from LiDAR sensors) is developed in this research, and semi-supervised learning is utilized to reduce the considerable requirement for fine annotations. The pairwise constraints

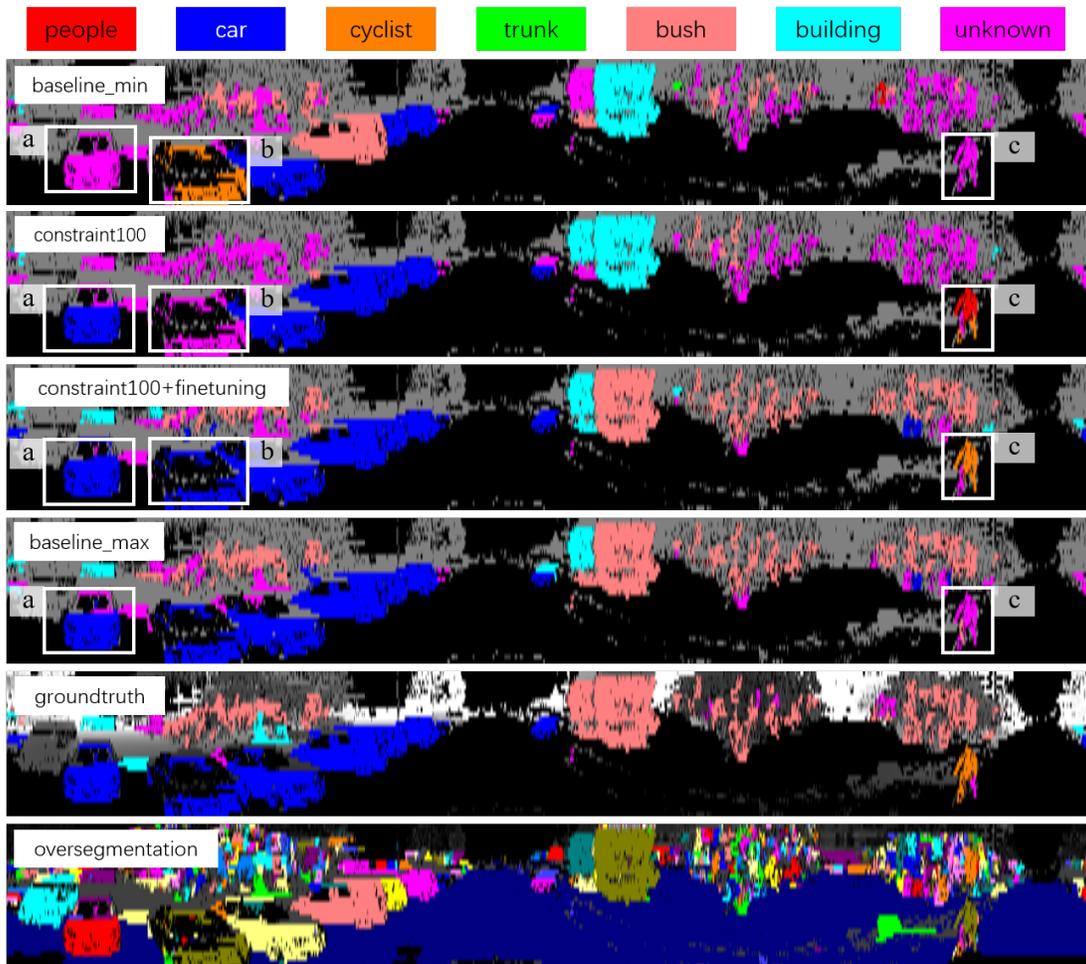


Fig. 12. The qualitative results on testing data. The a and b show the comparison on car and c shows that the cyclist is successfully classified after fine-tuning.

between adjacent frames are generated via inter-frame data association, and a loss function is designed to help the constraint data to obtain the same label. This method is examined extensively on a new dataset. The superior results indicate improvements in effectiveness and adaptability. Future work will address how to define the sample because including both foreground and background information in the sample can confuse the classifier. In addition, the introduction of new constraints will also be studied.

REFERENCES

- [1] C. Urmson, J. Anhalt, D. Bagnell, C. Baker, R. Bittner, M. Clark, J. Dolan, D. Duggins, T. Galatali, C. Geyer *et al.*, “Autonomous driving in urban environments: Boss and the urban challenge,” *Journal of Field Robotics*, vol. 25, no. 8, pp. 425–466, 2008.
- [2] F. Moosmann, O. Pink, and C. Stiller, “Segmentation of 3d lidar data in non-flat urban environments using a local convexity criterion,” in *IEEE Intelligent Vehicles Symposium*. IEEE, 2009, pp. 215–220.
- [3] B. Douillard, J. Underwood, N. Kuntz, V. Vlaskine, A. Quadros, P. Morton, and A. Frenkel, “On the segmentation of 3d lidar point clouds,” in *IEEE International Conference on Robotics and Automation*. IEEE, 2011, pp. 2798–2805.
- [4] D. Munoz, N. Vandapel, and M. Hebert, “Onboard contextual classification of 3-d point clouds with learned high-order markov random fields,” in *IEEE international conference on Robotics and Automation*. IEEE, 2009, pp. 4273–4280.
- [5] H. Zhao, Y. Liu, X. Zhu, Y. Zhao, and H. Zha, “Scene understanding in a large dynamic environment through a laser-based sensing,” in *IEEE International Conference on Robotics and Automation*. IEEE, 2010, pp. 127–133.
- [6] A. Garcia-Garcia, S. Orts-Escolano, S. Oprea, V. Villena-Martinez, and J. Garcia-Rodriguez, “A review on deep learning techniques applied to semantic segmentation,” *arXiv preprint arXiv:1704.06857*, 2017.
- [7] X. Zhu, “Semi-supervised learning literature survey,” *Computer Science, University of Wisconsin-Madison*, vol. 2, no. 3, p. 4, 2006.
- [8] T. Lange, M. H. Law, A. K. Jain, and J. M. Buhmann, “Learning with constrained and unlabelled data,” in *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 1. IEEE, 2005, pp. 731–738.
- [9] Z. Lu and T. K. Leen, “Semi-supervised learning with penalized probabilistic clustering,” in *Advances in Neural Information Processing Systems*, 2005, pp. 849–856.
- [10] M. Weinmann, B. Jutzi, and C. Mallet, “Semantic 3d scene interpretation: A framework combining optimal neighborhood size selection with relevant features,” *ISPRS Annals of Photogrammetry, Remote Sensing and Spatial Information Sciences*, vol. II-3, pp. 181–188, 2014.
- [11] D. Munoz, J. A. Bagnell, N. Vandapel, and M. Hebert, “Contextual classification with functional max-margin markov networks,” in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2009, pp. 975–982.
- [12] A. Serna, B. Marcotegui, F. Goulette, and J.-E. Deschaud, “Paris-rue-madame database: a 3d mobile laser scanner dataset for benchmarking urban detection, segmentation and classification methods,” in *International Conference on Pattern Recognition, Applications and Methods*, 2014.
- [13] T. Hackel, J. D. Wegner, and K. Schindler, “Fast semantic segmentation of 3d point clouds with strongly varying density,” *ISPRS Annals of*

- Photogrammetry, Remote Sensing and Spatial Information Sciences*, vol. III-3, pp. 177–184, 2016.
- [14] H. Hu, D. Munoz, J. A. Bagnell, and M. Hebert, “Efficient 3-d scene analysis from streaming data,” in *IEEE International Conference on Robotics and Automation*. IEEE, 2013, pp. 2297–2304.
- [15] J. Behley, V. Steinhage, and A. B. Cremers, “Performance of histogram descriptors for the classification of 3d laser range data in urban environments,” in *Robotics and Automation (ICRA), 2012 IEEE International Conference on*. IEEE, 2012, pp. 4391–4398.
- [16] Y. Lu and C. Rasmussen, “Simplified markov random fields for efficient semantic labeling of 3d point clouds,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2012, pp. 2690–2697.
- [17] F. Engelmann, T. Kontogianni, A. Hermans, and B. Leibe, “Exploring spatial context for 3d semantic segmentation of point clouds,” in *IEEE International Conference on Computer Vision Workshops*. IEEE, 2017, pp. 716–724.
- [18] I. Armeni, O. Sener, A. R. Zamir, H. Jiang, I. Brilakis, M. Fischer, and S. Savarese, “3d semantic parsing of large-scale indoor spaces,” in *International Conference on Computer Vision and Pattern Recognition*. IEEE, 2016.
- [19] A. Gaidon, Q. Wang, Y. Cabon, and E. Vig, “Virtual worlds as proxy for multi-object tracking analysis,” in *International Conference on Computer Vision and Pattern Recognition*. IEEE, 2016.
- [20] P. Tosteberg, “Semantic segmentation of point clouds using deep learning,” Master’s thesis, Linköping University, 2017.
- [21] L. University, “Virtual photo sets,” <http://www.hdrv.org/vps/>, accessed June 10, 2018.
- [22] G. L. O. a. B. Ayush Dewan, “Deep semantic classification for 3d lidar data,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2017, pp. 3544–3549.
- [23] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, “Vision meets robotics: The kitti dataset,” *International Journal of Robotics Research (IJRR)*, 2013.
- [24] T. Hackel, N. Savinov, L. Ladicky, J. D. Wegner, K. Schindler, and M. Pollefeys, “SEMANTIC3D.NET: A new large-scale point cloud classification benchmark,” *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, vol. IV-1-W1, pp. 91–98, 2017.
- [25] B. Wu, A. Wan, X. Yue, and K. Keutzer, “Squeezeseg: Convolutional neural nets with recurrent crf for real-time road-object segmentation from 3d lidar point cloud,” *arXiv preprint arXiv:1710.07368*, 2017.
- [26] F. Piewak, P. Pinggera, M. Schäfer, D. Peter, B. Schwarz, N. Schneider, D. Pfeiffer, M. Enzweiler, and M. Zöllner, “Boosting lidar-based semantic labeling by cross-modal training data generation,” *arXiv preprint arXiv:1804.09915*, 2018.
- [27] L. Caltagirone, S. Scheidegger, L. Svensson, and M. Wahde, “Fast lidar-based road detection using fully convolutional neural networks,” in *IEEE Intelligent Vehicles Symposium*. IEEE, 2017, pp. 1019–1024.
- [28] F. J. Lawin, M. Danelljan, P. Tosteberg, G. Bhat, F. S. Khan, and M. Felsberg, “Deep projective 3d semantic segmentation,” in *International Conference on Computer Analysis of Images and Patterns*. Springer, 2017, pp. 95–107.
- [29] L. P. Tchapmi, C. B. Choy, I. Armeni, J. Gwak, and S. Savarese, “Segcloud: Semantic segmentation of 3d point clouds,” *arXiv preprint arXiv:1710.07563*, 2017.
- [30] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, “Indoor segmentation and support inference from rgb-d images,” in *European Conference on Computer Vision*. Springer, 2012, pp. 746–760.
- [31] G. Riegler, A. O. Ulusoy, and A. Geiger, “Octnet: Learning deep 3d representations at high resolutions,” in *IEEE Conference on Computer Vision and Pattern Recognition*, vol. 3. IEEE, 2017, pp. 6620–6629.
- [32] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, “3d shapenets: A deep representation for volumetric shapes,” in *IEEE conference on computer vision and pattern recognition*. IEEE, 2015, pp. 1912–1920.
- [33] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, “Pointnet: Deep learning on point sets for 3d classification and segmentation,” in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2017, pp. 77–85.
- [34] L. Yi, V. G. Kim, D. Ceylan, I. Shen, M. Yan, H. Su, C. Lu, Q. Huang, A. Sheffer, L. Guibas *et al.*, “A scalable active framework for region annotation in 3d shape collections,” *ACM Transactions on Graphics (TOG)*, vol. 35, no. 6, p. 210, 2016.
- [35] I. Armeni, O. Sener, A. R. Zamir, H. Jiang, I. Brilakis, M. Fischer, and S. Savarese, “3d semantic parsing of large-scale indoor spaces,” in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2016, pp. 1534–1543.
- [36] L. Landrieu and M. Simonovsky, “Large-scale point cloud semantic segmentation with superpoint graphs,” *arXiv preprint arXiv:1711.09869*, 2017.
- [37] A. Bearman, O. Russakovsky, V. Ferrari, and L. Fei-Fei, “Whats the point: Semantic segmentation with point supervision,” in *European Conference on Computer Vision*. Springer, 2016, pp. 549–565.
- [38] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, “The pascal visual object classes (voc) challenge,” *International journal of computer vision*, vol. 88, no. 2, pp. 303–338, 2010.
- [39] R. Yan, J. Zhang, J. Yang, and A. G. Hauptmann, “A discriminative learning framework with pairwise constraints for video object classification,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 28, no. 4, pp. 578–593, 2006.
- [40] M. Bäuml, M. Tapaswi, and R. Stiefelhagen, “Semi-supervised learning with constraints for person identification in multimedia data,” in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2013, pp. 3602–3609.
- [41] S. Hong, H. Noh, and B. Han, “Decoupled deep neural network for semi-supervised semantic segmentation,” in *Advances in Neural Information Processing Systems*, 2015, pp. 1495–1503.
- [42] G. Papandreou, L.-C. Chen, K. P. Murphy, and A. L. Yuille, “Weakly-and semi-supervised learning of a deep convolutional network for semantic image segmentation,” in *IEEE International Conference on Computer Vision*. IEEE, 2015, pp. 1742–1750.
- [43] D. Lin, J. Dai, J. Jia, K. He, and J. Sun, “Scribblesup: Scribble-supervised convolutional networks for semantic segmentation,” in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2016, pp. 3159–3167.
- [44] T. Cour, B. Sapp, C. Jordan, and B. Taskar, “Learning from ambiguously labeled images,” in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2009, pp. 919–926.
- [45] D. Pathak, P. Krahenbuhl, and T. Darrell, “Constrained convolutional neural networks for weakly supervised segmentation,” in *IEEE International Conference on Computer Vision*. IEEE, 2015, pp. 1796–1804.
- [46] J. Xu, A. G. Schwing, and R. Urtasun, “Learning to segment under various forms of weak supervision,” in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2015, pp. 3781–3790.
- [47] C. Liu, J. Yuen, and A. Torralba, “Nonparametric scene parsing via label transfer,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, no. 12, pp. 2368–2382, 2011.
- [48] J. Dai, K. He, and J. Sun, “Boxsup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation,” in *IEEE International Conference on Computer Vision*. IEEE, 2015, pp. 1635–1643.