

Deep Active Learning for Efficient Training of a LiDAR 3D Object Detector

Di Feng^{1,2}, Xiao Wei^{1,3}, Lars Rosenbaum¹, Atsuto Maki³, Klaus Dietmayer²

Abstract—Training a deep object detector for autonomous driving requires a huge amount of labeled data. While recording data via on-board sensors such as camera or LiDAR is relatively easy, annotating data is very tedious and time-consuming, especially when dealing with 3D LiDAR points or radar data. Active learning has the potential to minimize human annotation efforts while maximizing the object detector’s performance. In this work, we propose an active learning method to train a LiDAR 3D object detector with the least amount of labeled training data necessary. The detector leverages 2D region proposals generated from the RGB images to reduce the search space of objects and speed up the learning process. Experiments show that our proposed method works under different uncertainty estimations and query functions, and can save up to 60% of the labeling efforts while reaching the same network performance.

Keywords—Deep neural network, active learning, uncertainty estimation, object detection, autonomous driving

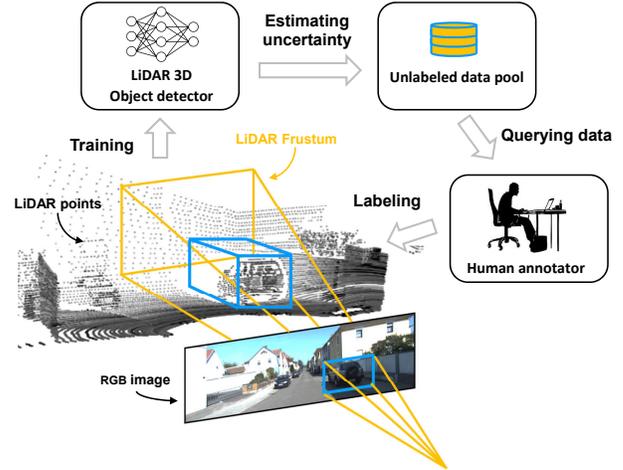


Fig. 1: Our proposed active learning method to efficiently train a LiDAR 3D object detector. The detector is based on 2D proposals from images, which serve as seeds to locate objects as frustums in LiDAR 3D space. We assume that there exists a large unlabeled data pool of LiDAR point clouds. The object detector iteratively uses predictive uncertainty to quantify the informativeness of each sample in the unlabeled data pool, queries the human annotator for the class label and 3D geometrical information of objects, and updates the training set with the newly-labeled data. We validate our method both with “perfect” image proposals provided by human annotators, or by an on-the-shelf pre-trained image detector with high recall rate.

I. INTRODUCTION

In recent years deep learning has set the benchmark for object detection tasks on many open datasets (e.g. KITTI [1], Cityscapes [2]), and has become the de-facto method for perception in autonomous driving. Despite its high performance, training a deep object detector usually requires a huge amount of labeled samples. The annotation process is tedious and time-consuming work, especially for 3D LiDAR points (as discussed in [3]), necessitating the development of methods to reduce labeling efforts. Furthermore, a common way to optimize a deep object detector is to feed all training samples into the network with random shuffling. However, the informativeness of each training sample differs, i.e. some are more informative and contribute more to the performance gain, while some others are less informative. A more efficient training strategy is to optimize the network with only the most informative samples. This is specifically helpful when adapting an object detector to new driving scenarios which are different from the previous training set, e.g. from highway to urban scenarios.

Active learning is a training strategy to reduce human annotation efforts while maximizing the performance of a machine learning model (usually in a supervised-learning

fashion) [4]. In active learning, a model iteratively evaluates the informativeness of unlabeled data, selects the most informative samples to be labeled by human annotators, and updates the training set with the newly-labeled data. Active learning has long been applied to Support Vector Machines (SVM) or Gaussian Processes (GP) [5]–[7], and has only recently been used in deep learning for classification of medical images [8] or hyperspectral images in remote sensing [9], 2D image detection [10], [11], road-scene image segmentation [12], and natural language processing [13].

In this work, we propose an active learning method to efficiently train a LiDAR 3D object detector for autonomous driving, as in Fig. 1. We assume that there exists a large unlabeled data pool of LiDAR point clouds, because it is relatively easy to collect and prepare LiDAR data using a test vehicle. We use the network’s predictive uncertainty to

¹ Robert Bosch GmbH, Corporate Research, Driver Assistance Systems and Automated Driving, 71272 Renningen, Germany.

² Institute of Measurement, Control and Microtechnology, Ulm University, 89081 Ulm, Germany.

³ School of Electrical Engineering and Computer Science, KTH Royal Institute of Technology, 100 44 Stockholm, Sweden.

quantify the informativeness of each sample in the unlabeled data pool, and assume that the network can iteratively query the human annotator for the class label and 3D geometrical information of objects. Furthermore, as it is much easier to do human labeling with 2D RGB images than 3D point clouds and a lot of pre-trained image detectors with high recall rate exist (e.g. Detectron [14]), we propose to leverage an image detector to provide 2D object proposals which serve as seeds to locate objects, so that the human annotator only needs to label LiDAR points within frustums (see Fig. 1). In this way, the 3D labeling efforts can be further reduced and the speed of learning process can be increased. In the experiments, we evaluate our method either by assuming a “perfect” image detector which provides accurate object proposals, or by an on-the-shelf pre-trained image detector. Results show our method outperforms the baseline method in both experimental settings.

Our **contributions** can be summarized as follows: (1) We propose a deep active learning method to significantly reduce the labeling efforts for training a 3D object detector using LiDAR points. To our knowledge, ours is the first attempt to introduce deep active learning for the 3D environment perception on autonomous driving. (2) Our method leverages the 2D object proposals from RGB images, which reduces the search space of objects of interests and speeds up the learning process. (3) We compare several approaches for quantifying uncertainties in the neural network, and study their efficiencies to query informative unlabeled data.

II. RELATED WORKS

In this section, we briefly summarize existing works on deep object detection for autonomous driving using LiDAR points as well as deep active learning.

A. Object Detection for Autonomous Driving using LiDAR points

Most driverless cars are equipped with multiple sensors, such as cameras and LiDARs. Therefore, many methods have been proposed for object detection using camera images [15]–[17], LiDAR point clouds [18]–[22], or the fusion of both to exploit their complementary properties [23]–[26].

State-of-the-art deep object detection networks follow two pipelines: the two-stage and the one-stage object detections. In the former pipeline, several object candidates called regions of interest (ROI) or region proposals (RP) are extracted from a scene. Then, these candidates are verified and refined in terms of classification scores and locations. For example, Asvadi *et al.* [27] cluster LiDAR points for on-ground obstacles using DBSCAN. These clusters are then fed into a ConvNet for 2D detection. Chen *et al.* [23] propose to generate 3D ROIs from the birds eye view LiDAR feature maps by a Region Proposal Network (RPN), and combine the regional features from the front view LiDAR feature

maps and RGB camera images for 3D vehicle detection. In the latter pipeline, single-stage and unified CNN models are used to directly map the input features to the detection outputs. Li *et al.* [19] and Yang *et al.* [28] employ the Fully Convolutional Network (FCN) on LiDAR point clouds to produce an objectness map and several bounding box maps. Caltagirone *et al.* [29] use a FCN for road detection. In this work, we follow the 2-stage object detection pipeline: 2D region proposals are provided by camera images, and the network detects objects using the LiDAR points in the corresponding frustum (Fig. 1), similar to [26].

B. Deep Active Learning

Active learning has a long history in the machine learning community (a comprehensive survey is provided by [4]), and has been introduced in deep neural networks in 2015 [8]. While many works exist in image classification [8], [9], [30] and segmentation problems [12], [31], [32], little attention has been paid to 2D image detection [10], [11]. Compared to these works, ours is the first attempt for active learning in 3D object detection problem.

There are numerous approaches to querying unlabeled data, such as variance reduction, query-by-committee, and expected model change [4]. Among them, the uncertainty-based approach suggests to use the predictive uncertainty to represent the data informativeness, and to query samples with the highest uncertainty. The effectiveness of this strategy is naturally dependent on obtaining reliable uncertainty estimates. Many recent methods have focused on obtaining such estimates in an efficient manner in deep neural networks. For example, Lakshminarayanan *et al.* [33] propose to use an ensemble of networks to predict uncertainty. Kendall *et al.* [34] decompose predictive uncertainty into model dependent (epistemic) and data dependent (aleatoric) uncertainties in Bayesian Neural Network. The former is obtained by Monte-Carlo dropout sampling [35], while the latter by predicting the noise in the input data. Guo *et al.* [36] use a simple calibration technique to improve the network’s probability output. Application of these uncertainties has also featured in many recent works. For example, Miller *et al.* [37] employ epistemic uncertainty for object detection in open-set scenarios. Feng *et al.* [38] evaluate the uncertainty estimation in a LiDAR 3D object detection network, and leverage the aleatoric uncertainty to significantly improve the network’s robustness against noisy data [22]. Ilg *et al.* [39] compare several uncertainty estimation methods in optical flow. In this work, we use Monte-Carlo dropout [35] and Deep Ensembles [33] to estimate uncertainties, and compare their efficiencies in query functions.

III. METHODOLOGY

In this work, we propose an active learning method to iteratively train a 3D LiDAR detector using the fewest

Algorithm 1 Active Learning for 3D LiDAR Object Detector

Input : $\mathcal{D}^u, \mathcal{D}^l, A$

```

1 Initialization:  $M \leftarrow \text{trainDetector}(\mathcal{D}^l)$ 
  while not StopCondition() do
2    $\mathcal{U}(\mathcal{D}^u) \leftarrow \text{uncertaintyEstimate}(\mathcal{D}^u)$ 
3    $X_u^* \leftarrow \text{dataQuery}(\mathcal{D}^u, \mathcal{U}(\mathcal{D}^u))$   $\triangleright$  A subset of unlabeled
     data
4    $Y^* \leftarrow \text{dataLabel}(X_u^*, A)$   $\triangleright$  Class label and object location
5    $\mathcal{D}^l \leftarrow \mathcal{D}^l \cup \{X_u^*, Y^*\}$   $\triangleright$  Add data to the training dataset
6    $\mathcal{D}^u \leftarrow \mathcal{D}^u \setminus X_u^*$   $\triangleright$  Delete data from the unlabeled dataset
7    $M \leftarrow \text{trainDetector}(\mathcal{D}^l)$   $\triangleright$  Update the network
8 end
  
```

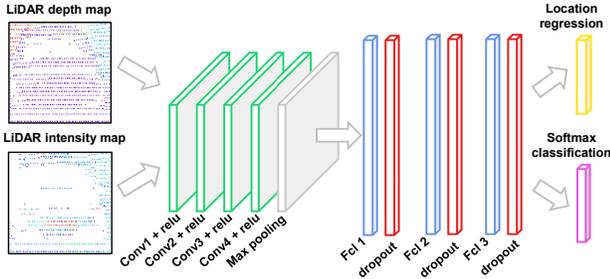


Fig. 2: Network architecture. The detector takes the LiDAR depth and intensity maps as input and outputs the objectness score and object location information (width, length, height, and depth).

number of labeled training samples, given a large unlabeled data pool.

Denote the large unlabeled data pool as \mathcal{D}^u , which consists of N^u i.i.d data samples $\mathcal{D}^u = \{\mathbf{x}_n^u\}_{n=1}^{N^u}$, and the labeled training dataset as \mathcal{D}^l , with N^l samples ($N^u \gg N^l$) and their labels $\mathcal{D}^l = \{\mathbf{x}_n^l, \mathbf{y}_n^l\}_{n=1}^{N^l}$. Also denote the human annotator as A and the object detection model as M . Our method is summarized in Alg. 1. The keys to the approach are the uncertainty estimation in neural networks, and the data query functions.

A. Process

To start the active learning process, the network is initialized with a small labeled dataset \mathcal{D}^l and trained in loop. After each training, the detector evaluates the informativeness of each sample in the unlabeled dataset \mathcal{D}^u by predicting the uncertainty (Step 2 in Alg. 1), selects the most informative samples via a query function (Step 3), and asks the human annotator for their class labels and object positions (Step 4). Afterwards, these samples are added to the labeled dataset, and the detector is re-trained. The process iterates until a stop condition is satisfied, e.g. the network’s performance converges for several iterations, or a desired performance is reached.

B. 3D Object Detector

1) *Inputs and Outputs:* As mentioned in the Introduction (Sec. I), our detector leverages the 2D region proposals in RGB images, which build frustums in the 3D space. We project those LiDAR points in frustums onto the front-view camera plane and build sparse LiDAR depth and intensity maps. Since these maps bring complementary LiDAR information (see Fig. 2), we concatenate them to build the network inputs. Note that we do not perform interpolation (e.g. Delaunay Triangulation [27] and Bilateral Filtering [40]) in order to avoid interpolation artifacts. The network outputs softmax classification scores \mathbf{s} , and object locations \mathbf{t} . We encode an object’s 3D position as the relative width $\hat{w} = \frac{w}{w_{max}}$, length $\hat{l} = \frac{l}{l_{max}}$, and height $\hat{h} = \frac{h}{h_{max}}$ of the bounding box, as well as the euclidean distance between our ego-vehicle and the object centroid $\hat{d} = \frac{d}{d_{max}}$, i.e. $\mathbf{t} = \{\hat{w}, \hat{l}, \hat{h}, \hat{d}\}$. We select $w_{max}, l_{max}, h_{max}, d_{max}$ based on the heuristics from the dataset.

2) *Network Architecture:* Our object detector is built on the ConvNet depicted in Fig. 2. It is composed of four convolutional layers (each with $32 \ 3 \times 3$ kernels and relu activation), a pooling layer, three fully connected layers (each with 256 hidden units), and three dropout layers. The dropout layers are used for stochastic regularization during training and uncertainty estimation during testing.

C. Uncertainty Estimation and Query Functions

1) *Uncertainty Estimation:* In this work, we use the predictive probability $p(\mathbf{y}|\mathbf{x})$ in classification to estimate uncertainty in our object detection network. For simplification, we denote \mathbf{x} as a data point and \mathbf{y} classification labels. A direct way to obtain predictive probability is by softmax output, i.e. $p(\mathbf{y}|\mathbf{x}) = \text{softmax}(\mathbf{x})$. However, as discussed in several works (e.g. [35], [30]), the softmax output may assign high probability to unseen data, resulting in over-confident predictions. Therefore, we also use two recent methods to obtain uncertainty estimates, namely, Monte-Carlo dropout (MC-dropout [35]) and Deep Ensembles ([33]).

MC-dropout [35] regards dropout regularization as approximate variational inference in the Bayesian Neural Network framework, and extracts predictive uncertainty by performing multiple feed-forward passes with dropout active during test time. More specifically, given a test point \mathbf{x} , the network performs T inferences with the same dropout rate as training, and averages the outputs to approximate the predictive probability:

$$p(\mathbf{y}|\mathbf{x}) \approx \frac{1}{T} \sum_{t=1}^T p(\mathbf{y}|\mathbf{x}, \mathbf{W}_t) = \frac{1}{T} \sum_{t=1}^T \text{softmax}_{(\mathbf{W}_t)}(\mathbf{x}), \quad (1)$$

with \mathbf{W}_t being network’s weights for the t^{th} inference.

Compared to MC-dropout, Deep Ensembles [33] estimates predictive uncertainty in a non-Bayesian way. It suggests

to train several networks with the same architecture but with random initialization, and average the networks’ outputs during testing. Let E be the number of ensembles and M_e a single network in the ensemble, similar to Eq. 1, we have:

$$p(\mathbf{y}|\mathbf{x}) \approx \frac{1}{E} \sum_{e=1}^E p(\mathbf{y}|\mathbf{x}, \mathbf{M}_e) = \frac{1}{E} \sum_{e=1}^E \text{softmax}_{(\mathbf{M}_e)}(\mathbf{x}). \quad (2)$$

2) *Query Functions*: Based on the above-mentioned methods to obtain the predictive probability, we can calculate the informativeness (or uncertainty) for each sample in the unlabeled data pool \mathcal{D}^u and use acquisition functions to query the most uncertain samples. A common way is to measure the Shannon Entropy (SE) [41] with:

$$\mathcal{H}[\mathbf{y}|\mathbf{x}] = - \sum_{c=1}^C p(y=c|\mathbf{x}) \log p(y=c|\mathbf{x}), \quad (3)$$

and query unlabeled samples with the highest Entropy values. y refers to a classification label, and C the number of classes.

Additionally, since both MC-dropout and Deep Ensembles provide samples from the predictive probability distribution (i.e. $p(\mathbf{y}|\mathbf{x}, \mathbf{W}_t)$ or $p(\mathbf{y}|\mathbf{x}, \mathbf{M}_e)$), we can use them to measure the Mutual Information (MI) [8] between the model weights and the class labels, and query unlabeled data with the highest MI. The mutual information using MC-dropout is calculated by:

$$\begin{aligned} \mathcal{I}[\mathbf{y}; \mathbf{W}] &= \mathcal{H}[\mathbf{y}|\mathbf{x}] - \mathbb{E}_{p(\mathbf{W}|\mathcal{D}^l)} \mathcal{H}[\mathbf{y}|\mathbf{x}, \mathbf{W}] \\ &\approx \mathcal{H}[\mathbf{y}|\mathbf{x}] + \frac{1}{T} \sum_{t=1}^T \sum_{c=1}^C p(y=c|\mathbf{x}, \mathbf{W}_t) \log p(y=c|\mathbf{x}, \mathbf{W}_t), \end{aligned} \quad (4)$$

where $p(\mathbf{W}|\mathcal{D}^l)$ indicates the posterior distribution of network weights \mathbf{W} given the training dataset calD^l . For Deep Ensembles, we only need to replace W with M . As discussed in our previous work [38], SE and MI capture different aspects of uncertainty: SE measures the output uncertainty (predictive uncertainty), whereas MI measures the model’s confidence in the data (epistemic uncertainty).

IV. EXPERIMENTAL RESULTS

A. Experimental Design

We evaluate our proposed method based on two experimental settings. In the first experiment, we study the active learning performance with different uncertainty estimation approaches and query functions. To avoid the influence from the RGB image detector, we assume a “perfect” image detector that provides only accurate object proposals. This is achieved by extracting objects using their ground-truth labels. Thus, we simplify the object detection problem in this setup to a classification and a location regression problem. In the second experiment, we use a pre-trained image detector to predict region proposals, which contain either object or background images.

Both experiments are conducted on the KITTI dataset [1]. The LiDAR depth and intensity maps are generated by projecting LiDAR points onto the image plane and then warped to 100×100 pixels, with additional 5 pixels padding to include context information, similar to [42]. To start the active learning process, the network is trained with some samples randomly selected from the training data and balanced over all classes (200 samples per class). The network is trained with Adam optimizer, together with dropout and l_2 regularization to prevent over-fitting. We set dropout rate to be 0.5, and the weight decay to be 10^{-4} . At each query step, 200 samples are selected from the remaining “unlabeled” training dataset, and the network is re-trained from scratch. At each query step, we calculate both the classification accuracy and the mean squared error (MSE) on the test data set, to track classification and localization performance, respectively. To have a fair comparison among different uncertainty estimation approaches, the MC-dropout methods are evaluated by performing the single forward pass on the test dataset without dropout. We also use the predictions from only one network in the ensemble for evaluation. For MC-Dropout, 20 forward passes are used during inference, while the ensemble consists of 5 classifiers. We fix the number of query steps to be 60, and repeat each experiment 3 times. Besides, we report the performance of a full-trained network.

B. Evaluation with a “Perfect” RGB Image Detector

1) *Setting*: We divide objects into five classes, namely, “Small Vehicle” (including “Car” and “Van” categories in KITTI), “Human” (including “Pedestrian”, “Person sitting”, and “Cyclist”), “Truck”, “Tram”, and “Misc”. We randomly divide the dataset into a training set with 31500 samples, a test set with 6000 samples and a val set with 3000 samples.

We compare the active learning strategy with the baseline method which randomly queries the unlabeled data points following the uniform distribution. We also study the behavior of active learning using five different query functions. “Softmax + Entropy”, “MC-dropout + Entropy” and “Ensembles + Entropy” use the same query function that maximizes Shannon Entropy. “MC-dropout + MI” and “Ensembles + MI” query the data by maximizing the mutual information.

2) *Active Learning Performance*: Results are shown in Fig. 3. We have three observations: (1) All active learning methods significantly outperform the baseline method for classification and localization tasks. They achieve higher recognition accuracy or lower mean squared error with the same number of training data as the baseline method. (2) With a relatively small amount of data (e.g. < 7000 samples), MI-based query functions (“MC-dropout + MI” and “Ensembles + MI”) consistently perform better than their Entropy counterparts (“MC-dropout + Entropy” and “Ensembles + Entropy”) in the localization task (Fig. 3(b)), whereas Entropy-based methods are more suitable for the

classification task (Fig. 3(a)). (3) Using MC-dropout and Deep Ensembles to estimate uncertainty results in slightly better active learning results compared to using a single softmax output (see “Softmax + Entropy”, “MC-dropout + Entropy” and “Ensembles + Entropy”).

Tab. I compares the number of labeled training samples required to train the detector so that it can reach a certain error level relative to the full-trained network. Denote $accu_{full}$ and $accu_m$ as the classification accuracy from the full-trained detector and the classification accuracy from a detector in the active learning process, respectively. Also denote mse_{full} and mse_m as the mean squared error for localization. We define the relative error for classification as $|accu_{full} - accu_m|$ and for localization as $|mse_m - mse_{full}|/mse_{full}$. We also calculate the percentage of labeling efforts saved by the active learning methods compared to the Baseline. As can be seen in the Table, our methods reach the relative error with significantly fewer training points, saving up to 60% training samples. In addition, using MC-dropout or Deep Ensembles to evaluate predictive uncertainty produces similar or better results than the single softmax approach, indicating that they can better represent the informativeness of unlabeled data.

3) *Understanding How Active Learning Works*: Our proposed method is based on a good uncertainty estimation. A better uncertainty score can better represent the data informativeness, leading to a better active learning performance. In this regard, we use the calibration plot and error curve to evaluate the quality of predictive uncertainty, similar to [30] and [39]. Furthermore, we investigate the class distribution of the queried samples.

Calibration Plot: A calibration plot is a quantile-quantile plot (QQ-plot) between the predictive probability of a model and the observed frequency of correct predictions. A well-calibrated uncertainty estimation should match the frequency of correct prediction, showing as a diagonal line. As an example, 60% of samples should be correctly classified as class c , when a well-calibrated network predicts them with probability output $p(y = c|x) = 60\%$. We compare the uncertainty estimations using single softmax (“Softmax + Entropy”), MC-dropout (“MC-dropout + Entropy”), and Deep Ensembles (“Ensemble + Entropy”) on the test dataset. The calibration plots at query steps 5, 15, 30, 55 are illustrated in the left four figures in Fig. 4, and the evolution of calibration error with query step in the right figure. The calibration error is calculated as the mean absolute deviation between the frequency and the diagonal line. The figures show that at the first few query steps all methods are “under-confident” with predictions, as their calibration plots are above the diagonal line. This indicates that networks are under-fitted with only a small number of training samples. As the query step increases, the predictive uncertainties from MC-dropout and Deep Ensembles become well-calibrated. However, the network using single softmax turns out to be over-fitted with data and produces over-confident predictions, resulting in a calibration plot under the diagonal line. The experiment show

that MC-dropout and Deep Ensembles produce more reliable uncertainty estimation and improve the learning performance more than the single softmax.

Error Curve: Another way to evaluate the quality of predictive uncertainty is via the error curve (or “sparsification plot” proposed in [39]). It is assumed that a well-estimated predictive uncertainty should correlate with the true error, and by gradually removing the predictions with high uncertainty, the average errors over the rest of the predictions will decrease. In our problem, we use the cross entropy loss to represent error. An exemplary error curve for “Ensemble + Entropy” at a specific query step is shown by Fig. 6(a). The benchmark is obtained by thresholding the predictions by their cross entropy losses (true errors). Note that for each uncertainty estimation method, we obtain a different benchmark. Therefore, we calculate the mean absolute deviation between the error curve and its benchmark, denoted as “Error sum”, to have a fair comparison on the quality of uncertainty estimation. Fig. 6(b) illustrates the evolution of the “Error sum” over the query steps for single softmax, MC-dropout and Deep Ensembles. Both MC-dropout and Deep Ensembles consistently outperform the single softmax method with smaller error sums.

Distribution of Sampled Objects To further understand how active learning outperforms the baseline method, we compare the class distribution of queried samples between “Ensemble + Entropy” and Baseline. This is calculated by taking the difference between the two at each step, normalized by the total number of samples from the corresponding classes in the unlabeled data pool. The results are shown in Fig. 5. The unlabeled data is highly class-imbalanced with ratios “Small Vehicle”= 78%, “Human”= 15.6%, “Truck”= 2.7%, “Tram”= 1.3%, “Misc”= 2.4%. However, our method naturally alleviates such problem by querying fewer samples from “Small Vehicle” and more from other classes. Note that this effect of balancing samples over classes is due to using the uncertainty estimation in the query function rather than an ad-hoc solution that explicitly over/undersample examples.

C. Evaluation with a Pre-trained RGB Image Detector

1) *Setting*: In this experiment, we evaluate the active learning method based on region proposals provided by a RGB image detector. To this end, we follow [23] to divide the KITTI dataset into a *train set* and a *val set*, and use the RGB image detector proposed by [26]. The *train set* is used to fine-tune the image detector which has been trained on COCO dataset [43], and the *val set* is used to evaluate our method. We consider to detect objects with the classes “Small vehicle” and “Human” (same to the previous experiment). A proposal is assigned positive when its 2D Intersection over Union (IoU) with the ground truth is larger than 0.5. Proposals with IoU smaller than 0.5 or from other object classes are marked as “Background”. The recall scores of the image detector is shown by Tab. II.

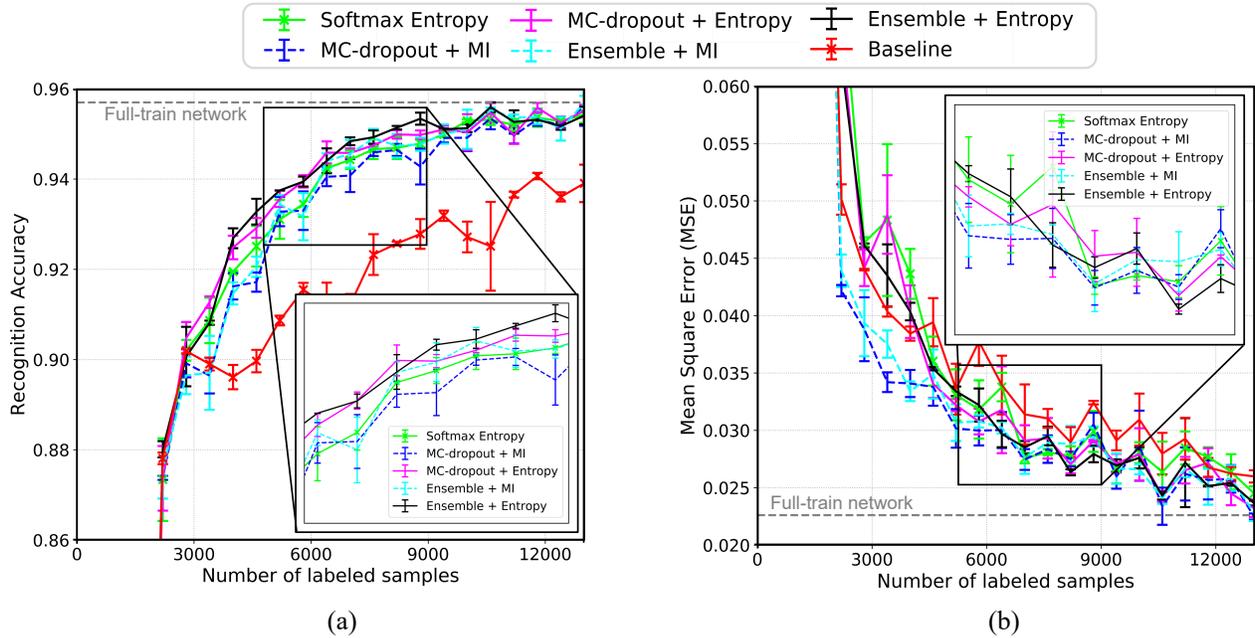


Fig. 3: Detection performance for the baseline and several active learning methods with different uncertainty estimations and query functions. The horizontal axis represents the increasing number of labeled training data samples. All networks are initialized with 1000 samples balanced over all classes. At each query step, 200 samples are queried from the unlabeled data pool. The vertical axis represents the detection performance for (a) classification and (b) object localization.

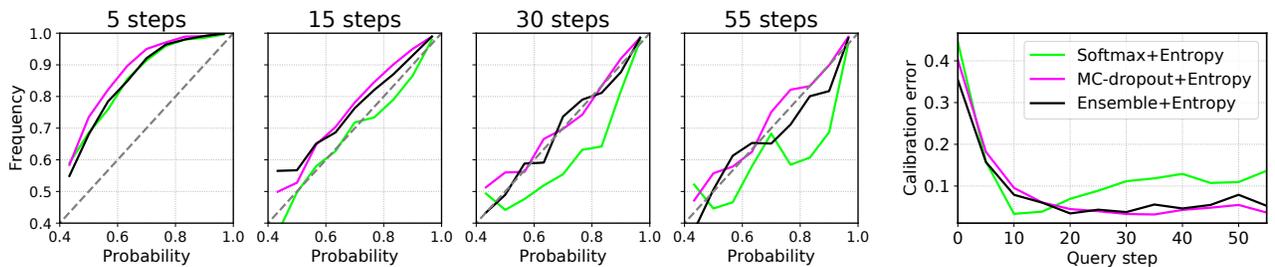


Fig. 4: A comparison of the calibration quality of predictive uncertainty from “Softmax + Entropy”, “MC-dropout + Entropy”, and “Ensemble + Entropy” averaged over three runs. To this end, we divide the probability values into several bins and calculate the frequency of correct predictions in each bin. The calibration plots at query step 5, 15, 30, and 55 are shown in the first four subplots. The diagonal line indicates perfect calibration, where the predictive uncertainty matches the observed frequency of correct predictions. The evolution of calibration error w.r.t. query step is shown in the right plot. A smaller error value indicates better uncertainty estimation.

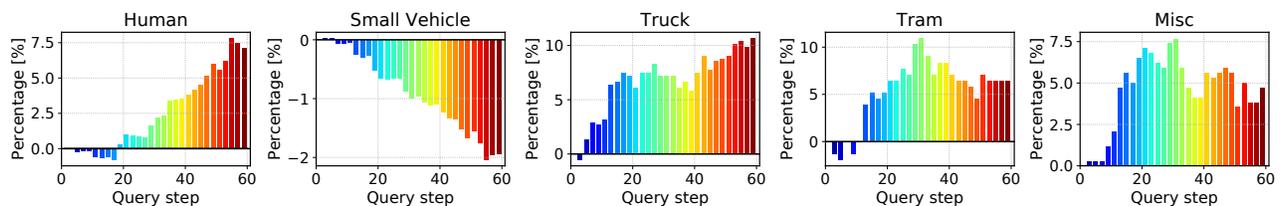


Fig. 5: Class distribution of queried samples for “Ensemble + Entropy” compared to the baseline method. The vertical axis represents the accumulative queried samples relative to the baseline, normalized by the number of samples of the corresponding classes in the unlabeled data pool. Compared to the baseline method, active learning queries fewer samples from “Small Vehicle” and more from the other classes. Note that since there are much more objects in “Small vehicle” than the other classes, a small percentage drop in “Small vehicle” means a much more percentage rise in other classes.

Relative error to the full-trained network	Recognition accuracy (Classification)				Mean Squared Error (Localization)			
	5%	4%	3%	2%	75%	30%	15%	5%
Baseline	5000	6800	9000	11400	4600	6800	10000	12200
Softmax + Entropy	3200(+36%)	3800(+44%)	4400(+51%)	6000(+47%)	3800(+17%)	6600(+3%)	8800(+12%)	10600(+13%)
MC-dropout + MI	3600(+28%)	3800(+44%)	4600(+49%)	6000(+47%)	2000(+57%)	5400(+21%)	8600(+14%)	10000(+18%)
MC-dropout + Entropy	2600(+48%)	3600(+47%)	4400(+51%)	5600(+51%)	3800(+17%)	6200(+9%)	8000(+20%)	10200(+16%)
Ensemble + MI	3400(+32%)	3800(+44%)	4400(+51%)	6000(+47%)	2200(+52%)	4400(+35%)	9000(+10%)	10200(+16%)
Ensemble + Entropy	3000(+40%)	3400(+50%)	3800(+58%)	4400(+61%)	3200(+30%)	4600(+32%)	7600(+24%)	10600(+13%)

TABLE I: The number of labeled training samples required to train the detector in order to reach a certain relative error to the full-trained network. The table also shows the percentage of labeling efforts saved by the active learning methods compared to the Baseline.

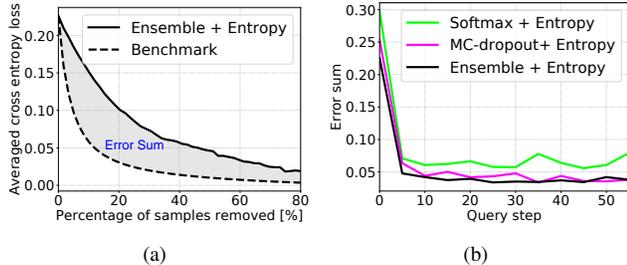


Fig. 6: (a): Error curve for “Ensemble + Entropy” at a specific query step. We use the cross entropy loss to represent the error. (b): The evolution of error sum w.r.t. query step. A smaller error sum indicates better uncertainty estimation. The plots are averaged over three runs.

	Small Vehicle	Human
Recall	0.917	0.862

TABLE II: Recall rate for the RGB image detector.

Based on image proposals, we build a training data pool with 17221 samples to train our active learning method. These samples are selected with their IoU being either larger than 0.5 (Positive) or below 0.2 (Background). The test dataset contains 6000 samples, including those with IoU ranging between 0.2 and 0.5. Note that ignoring samples with some IoU ranges is a common procedure when training object detectors, as discussed in [23], [42].

2) *Active Learning Performance*: We compare the active learning based on “Ensemble + Entropy” strategy with the baseline method. Results are shown in Fig. 7(a) and Fig. 7(b). Compared to the first experiment (Fig. 3), the network in this experiment results in a lower recognition accuracy, as this experiment is more challenging than the previous one. Despite of this, active learning consistently outperforms the baseline methods by reaching the same recognition accuracy or mean squared error with fewer labeled training samples.

V. DISCUSSION AND CONCLUSION

We presented a method that leverages active learning to efficiently train a LiDAR 3D object detector. The network predicts the object classification score and 3D geometric information based on 2D proposals on camera images. We conducted experiments using a “perfect” image detector to compare several ways of uncertainty estimation and

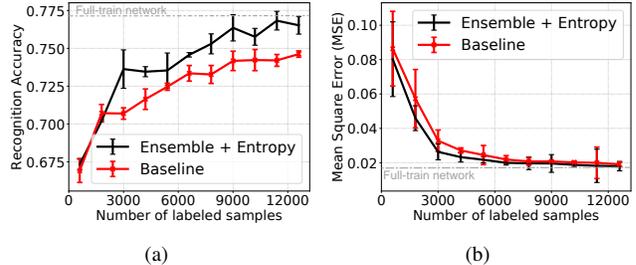


Fig. 7: A comparison of learning between active learning and baseline methods for (a) classification and (b) localization tasks. The active learning is built by “Ensemble + Entropy” strategy.

query functions. Results show that MC-dropout and Deep Ensembles provide more reliable predictive uncertainties compared to the single softmax output, and achieve better active learning performance. We also used a pre-trained image detector to predict image region proposals. In both experimental settings, our active learning method reaches the same detection performance with significantly fewer training samples compared to the baseline method, saving up to 60% labeling efforts.

We show that building query functions based on predictive uncertainty in classification is effective not only in improving recognition accuracy, but also in reducing the mean squared error for the localization task at the same time (e.g. Tab. I and Fig. 3). This indicates that by sharing weights in the hidden layers, the classification and localization are related to each other in the object detection network. It is an interesting future work to introduce location uncertainty into our active learning method. Furthermore, in this work the network is retrained from scratch after each query step. In applications where we want to adapt an object detector to new driving scenarios (as discussed in Sec. I), it is preferable to fine-tune the network with newly-labeled data. Employing our proposed active learning method in such a “life-long learning” scenario is an interesting future work.

One limitation of our method is that the performance of the LiDAR detector is highly dependent on region proposals from the image detection. Despite on-the-shelf image detectors already achieve high detection performance, the LiDAR detector can not handle false negatives in images. In order to guarantee highly qualified region proposals, we can incorporate the image detector into the active learning

loop, i.e. the region proposals are first provided by the image detector and then corrected by human annotators. We leave this as an interesting future work.

ACKNOWLEDGMENT

We thank William H. Beluch, Radek Mackowiak, and Christian H. Schuetz for the suggestions and fruitful discussions.

REFERENCES

- [1] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2012, pp. 3354–3361.
- [2] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [3] J. Lee, S. Walsh, A. Harakeh, and S. L. Waslander, "Leveraging pre-trained 3d object detection models for fast ground truth generation," in *2018 IEEE 21th International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2018.
- [4] B. Settles, "Active learning," *Synthesis Lectures on Artificial Intelligence and Machine Learning*, vol. 6, no. 1, pp. 1–114, 2012.
- [5] S. Tong and D. Koller, "Support vector machine active learning with applications to text classification," *Journal of machine learning research*, vol. 2, no. Nov, pp. 45–66, 2001.
- [6] A. Kapoor, K. Grauman, R. Urtasun, and T. Darrell, "Active learning with gaussian processes for object categorization," in *IEEE 11th International Conference on Computer Vision*. IEEE, 2007, pp. 1–8.
- [7] M. Kaboli, D. Feng, K. Yao, P. Lanillos, and G. Cheng, "A tactile-based framework for active object learning and discrimination using multimodal robotic skin," *IEEE Robotics and Automation Letters*, vol. 2, no. 4, pp. 2143–2150, 2017.
- [8] Y. Gal, R. Islam, and Z. Ghahramani, "Deep bayesian active learning with image data," in *International Conference on Machine Learning*, 2017, pp. 1183–1192.
- [9] J. M. Haut, M. E. Paoletti, J. Plaza, J. Li, and A. Plaza, "Active learning with convolutional neural networks for hyperspectral image classification using a new bayesian approach," *IEEE Transactions on Geoscience and Remote Sensing*, no. 99, pp. 1–22, 2018.
- [10] S. Roy, A. Unmesh, and V. P. Namboodiri, "Deep active learning for object detection," in *British Machine Vision Conference 2018, BMVC 2018, Northumbria University, Newcastle, UK, September 3-6, 2018*, 2018, p. 91.
- [11] C.-C. Kao, T.-Y. Lee, P. Sen, and M.-Y. Liu, "Localization-aware active learning for object detection," *arXiv preprint arXiv:1801.05124*, 2018.
- [12] R. Mackowiak, P. Lenz, O. Ghorri, F. Diego, O. Lange, and C. Rother, "Cereals-cost-effective region-based active learning for semantic segmentation," in *Proceedings of the British Machine Vision Conference (BMVC)*, 2018.
- [13] A. Siddhant and Z. C. Lipton, "Deep bayesian active learning for natural language processing: Results of a large-scale empirical study," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2018, pp. 2904–2909.
- [14] R. Girshick, I. Radosavovic, G. Gkioxari, P. Dollár, and K. He, "Detectron," <https://github.com/facebookresearch/detectron>, 2018.
- [15] X. Chen, K. Kundu, Z. Zhang, H. Ma, S. Fidler, and R. Urtasun, "Monocular 3d object detection for autonomous driving," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2147–2156.
- [16] A. Mousavian, D. Anguelov, J. Flynn, and J. Košecká, "3d bounding box estimation using deep learning and geometry," in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2017, pp. 5632–5640.
- [17] M. Teichmann, M. Weber, M. Zoellner, R. Cipolla, and R. Urtasun, "Multinet: Real-time joint semantic reasoning for autonomous driving," in *2018 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2018, pp. 1013–1020.
- [18] B. Li, "3d fully convolutional network for vehicle detection in point cloud," in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 1513–1518.
- [19] B. Li, T. Zhang, and T. Xia, "Vehicle detection from 3d lidar using fully convolutional network," *Robotics: Science and Systems*, 2016.
- [20] Y. Zhou and O. Tuzel, "Voxelnet: End-to-end learning for point cloud based 3d object detection," *arXiv preprint arXiv:1711.06396*, 2017.
- [21] M. Engelcke, D. Rao, D. Z. Wang, C. H. Tong, and I. Posner, "Vote3deep: Fast object detection in 3d point clouds using efficient convolutional neural networks," in *IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 1355–1361.
- [22] D. Feng, L. Rosenbaum, F. Timm, and K. Dietmayer, "Leveraging heteroscedastic aleatoric uncertainties for robust real-time lidar 3d object detection," *arXiv preprint arXiv:1809.05590*, 2018.
- [23] X. Chen, H. Ma, J. Wan, B. Li, and T. Xia, "Multi-view 3d object detection network for autonomous driving," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [24] J. Ku, M. Mozifian, J. Lee, A. Harakeh, and S. Waslander, "Joint 3d proposal generation and object detection from view aggregation," *arXiv preprint arXiv:1712.02294*, 2017.
- [25] D. Xu, D. Anguelov, and A. Jain, "Pointfusion: Deep sensor fusion for 3d bounding box estimation," *arXiv preprint arXiv:1711.10871*, 2017.
- [26] C. R. Qi, W. Liu, C. Wu, H. Su, and L. J. Guibas, "Frustum pointnets for 3d object detection from rgb-d data," *arXiv preprint arXiv:1711.08488*, 2017.
- [27] A. Asvadi, L. Garrote, C. Premebida, P. Peixoto, and U. J. Nunes, "Depthcn: vehicle detection using 3d-lidar and convnet," in *IEEE International Conference on Intelligent Transportation Systems (ITSC)*, 2017.
- [28] B. Yang, W. Luo, and R. Urtasun, "Pixor: Real-time 3d object detection from point clouds," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [29] L. Caltagirone, S. Scheidegger, L. Svensson, and M. Wahde, "Fast lidar-based road detection using fully convolutional neural networks," in *2017 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2017, pp. 1019–1024.
- [30] W. H. Beluch, T. Genewein, A. Nürnberger, and J. M. Köhler, "The power of ensembles for active learning in image classification," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [31] M. Gorriz, A. Carlier, E. Faure, and X. Giró i Nieto, "Localization-aware active learning for object detection," *arXiv preprint arXiv:1711.09168*, 2017.
- [32] K. Chitta, J. M. Alvarez, and A. Lesnikowski, "Large-scale visual active learning with deep probabilistic ensembles," 2018.
- [33] B. Lakshminarayanan, A. Pritzel, and C. Blundell, "Simple and scalable predictive uncertainty estimation using deep ensembles," in *Advances in Neural Information Processing Systems*, 2017, pp. 6402–6413.
- [34] A. Kendall and Y. Gal, "What uncertainties do we need in bayesian deep learning for computer vision?" in *Advances in Neural Information Processing Systems*, 2017, pp. 5580–5590.
- [35] Y. Gal, "Uncertainty in deep learning," Ph.D. dissertation, University of Cambridge, 2016.
- [36] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger, "On calibration of modern neural networks," in *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, 2017, pp. 1321–1330.
- [37] D. Miller, L. Nicholson, F. Dayoub, and N. Sünderhauf, "Dropout sampling for robust object detection in open-set conditions," in *2018 IEEE International Conference on Robotics and Automation, ICRA 2018, Brisbane, Australia, May 21-25, 2018*, 2018, pp. 1–7.
- [38] D. Feng, L. Rosenbaum, and K. Dietmayer, "Towards safe autonomous driving: Capture uncertainty in the deep neural network for lidar 3d vehicle detection," in *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, Nov 2018, pp. 3266–3273.
- [39] E. Ilg, O. Çiçek, S. Galesso, A. Klein, O. Makansi, F. Hutter, and T. Brox, "Uncertainty estimates and multi-hypotheses networks for optical flow," in *European Conference on Computer Vision (ECCV)*, 2018.
- [40] C. Premebida, L. Garrote, A. Asvadi, A. P. Ribeiro, and U. Nunes, "High-resolution lidar-based depth mapping using bilateral filter," in *2016 IEEE 19th International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2016, pp. 2469–2474.

- [41] C. E. Shannon, "A mathematical theory of communication," *ACM SIGMOBILE mobile computing and communications review*, vol. 5, no. 1, pp. 3–55, 2001.
- [42] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 580–587.
- [43] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755.