

## Map-aided annotation for pole base detection

Benjamin Missaoui, Maxime Noizet, Philippe Xu

### ▶ To cite this version:

Benjamin Missaoui, Maxime Noizet, Philippe Xu. Map-aided annotation for pole base detection. 35th IEEE Intelligent Vehicles Symposium (IV 2023), Jun 2023, Anchorage, AK, United States. 10.1109/IV55152.2023.10186774 . hal-04200337v3

### HAL Id: hal-04200337 https://hal.science/hal-04200337v3

Submitted on 1 Mar 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Map-aided annotation for pole base detection

Benjamin Missaoui<sup>1</sup>

Maxime Noizet<sup>1</sup>

Philippe Xu<sup>1</sup>

Abstract—For autonomous navigation, high definition maps are a widely used source of information. Pole-like features encoded in HD maps such as traffic signs, traffic lights or streetlights can be used as landmarks for localization. For this purpose, they first need to be detected by the vehicle using its embedded sensors. While geometric models can be used to process 3D point clouds retrieved by lidar sensors, modern imagebased approaches rely on deep neural network and therefore heavily depend on annotated training data. In this paper, a 2D HD map is used to automatically annotate pole-like features in images. In the absence of height information, the map features are represented as pole bases at the ground level. We show how an additional lidar sensor can be used to filter out occluded features and refine the ground projection. We also demonstrate how an object detector can be trained to detect a pole base. To evaluate our methodology, it is first validated with data manually annotated from semantic segmentation and then compared to our own automatically generated annotated data recorded in the city of Compiègne, France.

Erratum: In the original version [1], an error occurred in the accuracy evaluation of the different models studied and the evaluation method applied on the detection results was not clearly defined. In this revision, we offer a rectification to this segment, presenting updated results, especially in terms of Mean Absolute Errors (MAE).

#### I. INTRODUCTION

High definition (HD) maps provide a rich prior knowledge to automated navigation systems. Many different types of information are encoded in HD maps, often organized in layers. The road topology layer helps decision-making and motion planning while the features layer can be used for localization. For the latter, there are various kinds of HD maps.

One type of maps consists in low-level features such as point clouds or image-based key points/frames. These types of maps are often constructed from a SLAM perspective and their use may be restricted to the use of sensors similar to the ones employed during the mapping phase. Another constraint comes from the scalability of such maps as they are often heavy in terms of storage. However, one main advantage is that the mapped features being essentially geometric, the map can be constructed with a low level of human intervention.

Another type of maps consists in so-called vector map, which encodes features at a higher semantic level. The features are typically road infrastructure such as traffic signs, traffic lights, road markings or sometimes building footprints. The geometry of these features are often vectorized in 2D with geometric primitives such as points and lines. This enables these maps to be relatively light and can be easily scaled up

<sup>1</sup>The authors are with the Université de technologie de Compiègne, CNRS, Heudiasyc, France. name.surname@hds.utc.fr



Fig. 1: HD map of the city of Compiègne, France, containing generic pole-like elements such as traffic signs, bollards or street lights.

at a city level. Contrary to the previous case, these maps often require a higher amount of intervention from a human operator and could be harder to update automatically.

In this work, we consider only vector maps as they are sensor agnostic and could be provided in generic frameworks such as OpenStreetMap. The information encoded by road features such as traffic signs and traffic lights can be used for navigation but the knowledge of their spatial position can also be used for localization purposes. More generally, all georeferenced features can serve as landmarks to estimate the vehicle pose. To do so, the vehicle needs to be able to perceive these features from its own embedded sensors. Fig. 1 shows an HD map of the city of Compiègne, France, containing several types of point features such as traffic signs or traffic lights among others.

Lidars and cameras are the most common sensors used in the field of intelligent vehicles. Image-based methods for object detection or semantic segmentation have seen a tremendous improvement in past years thanks to deep neural networks. A key element for these methods to work is to have a large amount of annotated data.

We propose in this paper a method to make use of HD maps for automatic annotation of images. In particular, we show how the concept of pole base can be defined from road features and how they can be detected from an object detection point-of-view using state-of-the-art deep learning methods. The contributions of this paper are as follows:

- definition of a pole base class from HD map pole-like features
- projection, refinement and filtering of image annotations by the use of a lidar sensor
- · bounding box annotations for object-based detection
- experimental validation using both manually annotated data from semantic segmentation datasets and from real data and an HD map in the city of Compiègne, France

The rest of the paper is organized as follows. First, some related works are presented in section II. In section III, we introduce how HD map features can be projected onto images for automatic annotations and why additional refinement and filtering using a lidar are needed. To compare with manually annotated data, we show in section IV how semantic segmentation data can be used to achieve similar annotations. Next, in section V, we introduce the use of a bounding-box based object detector for pointwise detection. Finally, experimental results are detailed in section VI, first using data from semantic segmentation datasets and then from real unannotated data along with an HD map acquired in the city of Compiègne.

#### II. RELATED WORKS

The use of pole-like features in localization context is fairly common [2]-[4]. The most common sensors used in intelligent vehicles context are cameras and lidars. For lidars, geometric assumptions can be used in order to build polelike object detectors [5]-[7]. The lidar point intensity can also be used to detect reflective traffic signs [8]. Contrary to lidars, modern camera-based detectors mostly rely on machine learning techniques especially through the use of deep neural networks. Therefore, they heavily rely on available annotated data. The detection of road features such as traffic signs or traffic lights can be made using object detection methods requiring bounding box annotations. Many generic object detectors are trained from datasets such as the Microsoft COCO [9] or more specialized one such as KITTI [10]. The detection of polelike features, however, such as street lights is often considered from semantic segmentation point-of-view. This requires pixellevel annotations which are very costly. Several large scale semantic segmentation datasets are nevertheless available such SemanticKITTI [11], Cityspaces [12] or BDD100K [13].

In order to be less reliant on manually annotated data, Dong et al. [14] proposed to use range images constructed from a lidar in order to detect pole-like features which then served as pseudo-labels to train a deep neural network. Sun et al. [15] used an HD map to generate image-level labels, e.g., number of lanes in an image. Lee et al. [16] proposed a semi-automatic traffic landmark annotation using 3D road features from an HD map. Their goal was to provide initial annotations to accelerate a human annotation process. In this paper, we aim at using 2D HD maps to automatically annotate pole-like features in images while using an object detection approach.

#### III. HD MAP FEATURES PROJECTION AND FILTERING

Throughout this paper, the map features are considered to be encoded as georeferenced two-dimensional points without height information. In order to apply geometric transformation to these point features, their geodetic coordinates, i.e., longitude and latitude, are first converted into an East-North-Up (ENU) coordinates frame with its origin being in the vicinity of the driving sequences. Note that the "up" coordinate can be ignored as the altitude is set arbitrarily. A map feature is therefore represented as a 2D point <sup>M</sup>P with coordinates  $\begin{bmatrix} Mx, My \end{bmatrix}^{T}$ expressed in the map frame denoted by the superscript M. The map is therefore a set of n 2D points <sup>M</sup> $\mathcal{M} = { {}^{M}P_{i} }_{i=1,...,n}$ .

#### A. Projection of pole base onto image frames

In order to project the map features onto the image frame, they are first expressed in the vehicle frame which is defined as the body frame of the vehicle IMU placed above the center of the rear axle. A map feature point  ${}^{V}P$  in the vehicle frame is related to  ${}^{M}P$  as follows

$$\begin{bmatrix} \mathbf{v}_x \\ \mathbf{v}_y \end{bmatrix} = \begin{bmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{bmatrix} \begin{bmatrix} \mathbf{M}_x - x \\ \mathbf{M}_y - y \end{bmatrix},\tag{1}$$

where  $[x, y]^{\top}$  is the coordinates vector of the vehicle and  $\theta$  its heading. Note that the accuracy of the coordinates of  $^{\mathrm{V}}P$  gets more sensitive with respect to the heading estimate  $\theta$  when the map features  $^{\mathrm{M}}P$  gets further away from the vehicle position.

Because there is no height information available, it is not possible to project the map features directly onto an image. Instead, we decided to consider all the map features at ground level. That is to say, regardless of the true height of the map feature, we will only consider its projection onto the ground. For pole-like features such as bollards or street lights, it will correspond to their base, i.e., the contact point with the ground. For traffic signs, most of them are actually fixed on top of poles although these poles may not be mapped explicitly. Therefore, all the map features will be considered as represented in the image frame by their respective pole base point.

By assuming that the x-y plane of the vehicle frame is parallel to the ground and at fixed known height h, a map feature can be set at the ground level by setting

$$^{\mathsf{V}}z := -h. \tag{2}$$

These points can then be projected onto the camera image by making use of the camera intrinsic calibration parameters as well as its extrinsic ones with respect to the vehicle frame.

Fig. 2 illustrates the projection of some map features onto an image. The map feature points are supposed to represent the ground level projection of the map features. The green point on the right side pictures is a well-projected map feature. The red ones represent ill-projected map features due to the parallel ground hypothesis being wrong. The black points depict projected map features that are actually not visible from the camera due to occlusions. The yellow points are distant features that are both occluded but also wrongly projected as they are all the more sensitive to the ground plane estimation. Finally, the blue circle shows an example of an unmapped feature. This shows that the projection of map features onto an image is not straightforward and needs further filtering and refinement.



Fig. 2: Naive projection of map features onto an image frame. Only the green point on the right-hand side can be considered as a correct annotation.

#### B. Lidar-based refinement and filtering

In order to deal with the issues pictured in Fig. 2, we propose to use an additional lidar sensor to refine and filter the map feature projection. The lidar is used for two purposes. The first is to better estimate the ground surface on which to project the map features and the second is to assess whether or not a feature is visible in the image frame.

For this purpose, we make use of a Hesai Pandora sensor, which provides a 3D lidar sensor placed on top of five cameras: four grayscale wide-angle cameras to cover a  $360^{\circ}$  field-ofview and an additional front color camera. This setup allows to have almost the same field-of-view for both the lidar and the cameras. That is to say that if a feature is occluded from the camera perspective, it is also the case from the lidar one.

Although the two modalities share the same field-of-view, the lidar sensor is limited in terms of range. In our case, we consider a threshold of 150 meters above which we assume that the lidar cannot be used. Therefore, a first step was to remove the map features outside a radius of 150 meters around the vehicle. Note that these features could nevertheless still be visible in the image frame. This first filtering allows to remove the yellow crosses pictured in Fig. 2.

the yellow crosses pictured in Fig. 2. Next, the 3D point cloud  $C = \left\{ p_i = [x_i, y_i, z_i]^\top \right\}_{i=1,...,m}$  from the lidar is used to have a better estimate of the ground

from the lidar is used to have a better estimate of the ground surface. In this work, we use the Patchwork++ ground segmentation method proposed in [17]. This algorithm separates all the lidar points into two groups, a ground points one  $\mathcal{G}$  and a non-ground one  $\overline{\mathcal{G}}$ . The group composed of ground points is used to estimate the height of the map features. To do so, the map feature point  $^{V}P$  in the vehicle frame is transformed into a point  $^{L}P$  in the lidar frame using a similar equation as in (1) with the extrinsic calibration parameters encoding the relative position of the lidar w.r.t. the vehicle frame. Similarly to the vehicle localization, the calibration quality of the lidar, especially in terms of angle, has a higher impact of map features far away from the vehicle. Once a map feature has been projected in the lidar frame from a top down 2D space, we use a nearest neighbor approach to select the closest lidar



Fig. 3: Projection of map features onto an image frame using Patchwork++.

point, categorized as belonging to the ground, and use its z-coordinate as a height estimate instead of the one from (2):

<sup>L</sup>
$$z := z_k$$
, with  $k = \underset{i, p_i \in \mathcal{G}}{\operatorname{arg\,min}} \left\| {}^{\mathrm{L}}P - p_i \right\|_{2D}$ , (3)

where the  $\|\cdot\|_{2D}$  operator is the Euclidean distance but using only the *x*-*y* coordinates.

Fig. 3 shows the lidar point cloud projected onto the image. The green dots correspond to point segmented as ground while the others are in red. Using this new height estimation, the red crosses on the left side in Fig. 2 are now correctly projected at the ground level in Fig. 3.

The final step is to remove the map features that are occluded by some obstacles. To do so, we first project all the lidar points onto the image as well as the map features. For each map feature, we estimate the average depth using the lidar points that surround it in the image frame within a given radius. This depth is then compared with the true distance separating the map feature and the camera frame. The idea is that if a map feature is occluded by the obstacle, such as the bus in Fig. 3, then the estimated depth will be shorter than the true distance. By setting a threshold on the distance difference, the occluded features can be removed. The black crosses in Fig. 3 are correctly classified as occluded and can be filtered.

These different steps enable to have a cleaner projection of the map features. However, we can see in Fig. 3 that some errors still remain. Because the ground segmentation can also be erroneous some wrong projection can still occur. We can also see the projection of the lidar onto the image is not perfect despite using the calibration parameters provided by the sensor manufacturer. And finally, the unmapped features remain unannotated.

#### IV. POINT ANNOTATION FROM SEMANTIC SEGMENTATION

The automatic annotation presented in the previous section is still imperfect due to various noises from localization, calibration, map errors or imperfect ground estimation. In order to assess the performance of a deep learning method for the detection of pole base features, we also study the use of manually labeled data.



(a) Point annotation(b) Semantic segmentationFig. 4: Point annotations from semantic segmentation.

To our knowledge, there exist no dataset with such annotations. Fortunately, pole-like features can be found in datasets dedicated to semantic segmentation. In this study, we consider the BDD100K dataset [13]. In this dataset, we merge the three categories "pole", "traffic sign" and "traffic light" to represent our map features. The annotations are provided at the pixel level, while in our case, we wish to mimic the point-wise annotations retrieved from an HD map. Therefore, for each cluster of annotated map feature pixels, we need to compute the position of the pole base.

Like in the previous section, the base of a pole may not be visible due to occlusions. Thanks to the pixel-wise annotations, we can define a pole base as the lower part of a map feature cluster that lies on top of a ground pixels. The ground is defined by merging the following classes: "road", "sidewalk" and "terrain". In addition to filtering out occluded poles, we also added a minimal width to the pixel clusters so as to filter out far away features. Fig. 4 illustrates some results on the BDD100K datasets. The crosses represent the pole bases lying on the ground.

#### V. POLE BASE DETECTOR

In this study, we propose to formalize the detection of the pole base as an object detection problem. Therefore, bounding boxes are used to represent the pole base. At first glance, using a box to encode a point may not be relevant, especially for pole-like features that have thin appearances. In our case, the boxes are not "bounding" an object as in classical object detection approaches. Here, the boxes are used to encode the context surrounding the pole base while the centers of the boxes represent the points we aim at detecting. The size of the boxes is a meta-parameter studied in the experimental process. Using a bounding box formalism also allows to make use of highly efficient object detectors such as YOLO. In this paper, the YOLOv7 [18] algorithm is used.

To evaluate the performance of the detector, we consider two types of metrics. First, the traditional object detection metrics are used (mAP, precision and recall), in order to evaluate the ability of the model to predict the right boxes. Note that we will be using the notation mAP 0.5:0.95 throughout this article, which refers to the mAP at different IoU thresholds, a metric commonly used in object detection as less saturated. Then, to have a finer evaluation of the point-wise accuracy of the box centers, we compute the horizontal distance with respect to the point annotations. This choice is motivated by the fact that within a monocular-camera localization context, the horizontal

#### TABLE I: BDD100K dataset information

	Number of images	Number of poles
Training	7628	9357
Validation	372	909

coordinate of the detection points in the image frame can be used for bearing-only localization. Therefore, the accuracy along the horizontal axis becomes more important than along the vertical one. Nevertheless, a Euclidean 2D distance could also be used instead to compute the point-wise accuracy.

#### VI. EXPERIMENTAL RESULTS

Our experiments were divided into two parts. First, we validated the YOLO-based pole base detector using data from the BDD100K dataset. These two cases represented the ideal context where the pole base features could be assumed to be perfectly annotated manually. Then in a second step, we used an HD map to automatically annotate images. This was done using an experimental Renault ZOE car equipped with a Pandora sensor and driving in the city of Compiègne, France.

#### A. BDD100K

Before trying to learn a pole base detector using our HD map aided automatic annotations, that is known to be partially imperfect, we first validated it using data from the BDD100K dataset. The annotation procedure described in Sec. IV was applied to all the images from each dataset. In addition, for the validation set, the annotations were manually reviewed in order to remove some remaining erroneous annotations generated from the semantic segmentation ground truth. Table I gives the number of training and validation examples for each dataset.

As stated in the previous section, we used an object detection approach with bounding boxes. The pole base features were encoded by a box centered on it and with fixed size. The choice of the box width and height then arose. We decided to experiment on BDD100K with squared boxes of various sizes ranging from 60 to 400 pixels. Employing the official implementation of YOLOv7 with default hyperparameters, we set the IoU threshold to 0.5 during the non-maximum suppression step for evaluation purposes. The model was initialized with the weights from MS COCO. The training took 18h roughly on a single Tesla V100 32G GPU for 300 epochs. Most of the training converged after 100 epochs.

The predictions are accompanied by a confidence score. The lower this score, the less likely it is that an object is actually present at the location predicted by the bounding box. We can reject predictions with confidence score below a given value to remove potential wrong detections and improve the precision at the potential expense of recall.

We chose arbitrarily to apply a confidence threshold of 0.25 as it is unlikely that these low-confidence predictions genuinely correspond to objects. Besides, in our specific application context, maintaining a satisfactory precision is crucial.

TABLE II: Detection metrics and MAE after 100 epochs of training with different box sizes on the validation set of BDD100K. Precision and recall obtained are given using the confidence threshold 0.25.

Box size	mAP 0.5:0.95	Precision	Recall	MAE (px)
60x60	42.5	71.8%	58.7%	2.28
80x80	44.8	74.2%	55.0%	2.87
100x100	54.7	67.6%	<b>68.1</b> %	3.18
150x150	57.6	76.5%	61.8%	4.23
200x200	58.2	<b>76.7</b> %	61.9%	5.29
250x250	58.5	66.1%	65.6%	7.90
300x300	55.6	71.9%	63.7%	9.73
350x350	55.2	67.7%	61.2%	12.5
400x400	60.4	73.2%	62.8%	15.8

The detection metrics obtained with this threshold for all the models are detailed in Table II. When the boxes were too small, typically less than 100 pixels, the model exhibited poorer performance than other models, especially in terms of mAP. The performance of the detector gets better and stable starting from 100 px.

One plausible explanation for this could be that small boxes do not contain enough local context. When gradually increasing their size, the boxes aggregate more and more information about the surroundings of the pole until these bits of information become enough for the model to recognize and discriminate the pole. In the context of BDD100K, the option to increase box size emerges as a viable strategy to enhance performance, particularly since poles exhibit sufficient spacing, mitigating the risk of potential missdetections. When the box size exceeds 100px, it appears to contribute to improved precision and mAP at the expense of diminished recall. However, the detection performance can vary a lot between two box sizes.

We studied the localization accuracy of the coordinates of the center of the boxes w.r.t. the pole base annotations. As stated previously, we only focused on the horizontal accuracy along the x-coordinate in the image frame. Using all the true positive predicted bounding boxes, we computed the Mean Absolute Error (MAE) between the predicted xcoordinate and the labeled one. As visible in Table II, enlarging the box size introduces more uncertainty in the pole base position. Significantly increasing box size can enhance overall performance, as shown by the 400x400 model ranking as the best in mean Average Precision (mAP). But, it is imperative to note a substantial increase of MAE and this suggests a tradeoff.

In this study, we chose to focus on a detector with high precision in detection, prioritizing a low false positive rate, even at the cost of accepting a reasonable reduction in recall and localization accuracy. We consequently opted for the 200x200 model as it achieved the highest precision and ranked third in terms of mAP 0.5:0.95. It is important to note that another confidence threshold could have led to another decision.

Precision-Recall curves of all models are visible in Figure 5.



Fig. 5: Precision-Recall curves obtained after 100 epochs of training with different box sizes on the validation set of BDD100K. The points corresponding to values obtained in Table II using a confidence threshold of 0.25 are visible using crosses.

It shows that the choice of the bounding box size depends in reality deeply on the precision-recall trade-off we want to apply. The curves confirm that 60x60 and 80x80 on BDD100K seems insufficient to obtain sufficient recall when aiming a high precision. A minimal size of 100x100 is better for any trade-off and for a small increase of uncertainty. Besides due to high MAE, 350x350 and 400x400 box sizes should be avoided, leading to a potential choice between 100x100 and 300x300 depending on the chosen trade-off.

As pictured in the figure, the arbitrary threshold applied previously corresponds to different precision-recall trade-offs depending on the model used. That is why, in the best scenario a precision or recall objective should be defined, then a box size and the corresponding threshold should be chosen to maximize the other metric while keeping in mind the localization accuracy objective. The box size consequently become a new hyperparameter to reach the best detection performance.

#### B. City of Compiègne

1) Data collection: An experimental Renault ZOE vehicle was equipped with the following sensors:

- NovAtel SPAN-CPT GNSS/IMU with post-processed PPK computations for centimeter-level accuracy localization (50 Hz).
- Hesai Pandora 40-layer LiDAR with a horizontal angular resolution of 0.2° and integrating five synchronized cameras (10 Hz) (a front color camera, and four wide-angle mono cameras). For this work, only the color camera was used.

The extrinsic calibration between the Pandora lidar and the SPAN was obtained using a high-accuracy FARO Vantage laser tracker. For intrinsic calibration of cameras and the lidar,



Fig. 6: Experimental Renault ZOE vehicle equipped with a NovAtel SPAN-CPT GNSS/IMU and a Hesai Pandora sensor.

TABLE III: Heudiasyc dataset information

	Number of images	Number of poles
Train	15724	6724
Validation	605	2148

factory settings were used. Fig. 6 shows the experimental platform along with the different working frames.

The recorded sequence illustrated in Fig. 1 covered the entire city of Compiègne, France, with diverse road contexts. It lasted 40 minutes for a cumulative distance of 13.5 km. This dataset was recorded within the context of the ERASMO [19] European project. The HD map used in this experiment was acquired in 2019. In the meantime, some new poles were added and others were removed. Additionally to annotation errors due to localization, calibration or projection, some other errors were also due to erroneous HD map data.

2) Training sets: In the driving sequence illustrated in Fig. 1, the blue part of the trajectory was driven in an area where the environment was not mapped. Therefore, we used this part of the sequence as the validation set and manually annotated the images in order not to bias the evaluation part. For this part, only one image per second was kept and the images were further subsampled avoid having too many similar examples. The rest of the sequence was used as our training set and the HD map was used to automatically annotate the images.

The sets obtained are described in Table III. It is important to note that although the number of images seems larger compared to BDD100K, there is a lack of variability as they correspond to a unique driving sequence. Similarly, the 6724 annotated poles actually correspond to a subset of 1600 unique poles in the HD map.

*3) Map-aided annotation:* To annotate our images using our HD map and the lidar, we set the parameters of the filtering described in Sec. III to the values detailed in Table IV.

Fig. 8 illustrates examples of map-aided annotated images from our dataset in different situations. The first row corresponds to a naive projection of the HD map features with a

Maximum distance between a pole and the vehicle	150m
Image search radius for pole lidar-based filtering	20px
Depth difference for pole lidar-based filtering	5m

TABLE V: Detection metrics after 100 epochs of training on the validation set of Compiègne using a 200x200 box size. Precision and recall are given for a confidence threshold of 0.25.

Model	mAP 0.5:0.95	Prec.	Rec.
BDD100K	34.9%	85.6%	28.9%
Compiègne	31.6%	63.9%	50.1%

simple filtering of far away features. The second row adds a lidar-based filtering to remove the occluded features. The last row adds an additional ground segmentation to refine the height estimation. In average, adding each of these additional filtering improve the annotations but occasionally some noises can also be introduced due to wrongly segmented ground or bad depth estimation. In addition, the unmapped features shown in Fig. 8e can not be dealt with our current approach.

Using these automatically annotated training data, we trained a model in the same way as we did on the BDD100K dataset. Table V shows the detection performances on the Compiègne validation set comparing the model trained on BDD100K to the one trained in Compiègne using a 200x200 box size and the same confidence threshold 0.25 used previously.

The precision-recall curves for BDD100K and Compiègne models, utilizing 200x200 and 250x250 box sizes, are depicted in Figure 7. As suggested by Table V, the BDD100K models exhibit distinct performance on the Compiègne validation set. For any precision objective above 0.25, the 250x250 model outperforms the other, illustrating the dependence of the box size choice on the studied validation set. Even for the Compiègne models, the 250x250 box size appears to be a preferable choice across various precision-recall trade-offs.

Fig. 9 illustrates a detection example using both 250x250 models on a Compiègne validation image with a confidence threshold set to 0.25.

We observe notable differences in the performance of the BDD100K model compared to the results presented in the previous section, characterized by high precision and poor recall. The lower mAP attained by this model suggests a distinct precision-recall curve compared to the previous observations. The optimal box size for this arbitrary confidence threshold may vary from the one used on the BDD100K validation set and even across different precision-recall tradeoffs. In contrast, the Compiègne model demonstrates satisfactory performance at this arbitrary threshold, showcasing the feasibility of training a pole base detector using automatic annotations generated from a map. However, the relatively low mAP indicates that performance may fluctuate when adjusting the confidence threshold or IoU threshold.



Fig. 7: Precision-Recall curves obtained after 100 epochs of training with BDD100K and Compiègne models with 200x200 and 250x250 box sizes. The points corresponding to precision-recall values obtained using a confidence threshold of 0.25 are visible using crosses.

The main reason is likely to come from the imperfect annotations despite the lidar-based refinement and filtering. A second reason may come from the low variability of the training data which was extracted from a single driving sequence. However, the automatic annotation framework introduced in this work enables to easily add supplementary training data from new driving sequences without additional cost in terms of data labeling.

#### VII. CONCLUSION

In this paper, we introduced a framework to use an HD map to automatically annotate images. The map features were encoded as pole bases and projected onto images while a lidar was used to filter occluded features and to better estimate the ground surface. Although a pole base was annotated by a single point, we demonstrated that it could be detected using a traditional object detector using bounding boxes. This approach was first validated from annotations generated from semantic segmentation on the BDD100K. The method was then tested using annotations generated from an HD map of the city of Compiègne.

This preliminary work opens room to further improvements. The lidar-based filtering and refinement could be improved by tuning better the different parameters, but also with more sophisticated geometric consideration. The training data could be subsampled in a same way as the validation data to avoid training on too similar images with potentially wrong automated annotations leading to a loss of performance. The automated annotations, which are still imperfect, can also serve as an initial annotation either for a human operator or as an input for a second stage refinement. The ability to have a low-cost annotation pipeline could be exploited in many contexts such as generating training data in different driving weather or lighting conditions without having the need to label the data from scratch. The HD map used in this work also contained other features such as lane markings. This approach could also be extended to annotate data to train a lane marking detector. Finally, the use of the proposed pole base detector will be studied within a localization context in future work.

#### ACKNOWLEDGMENT

This work has been funded by the European project ERASMO [19] (GSA/GRANT/03/2018) in the framework of the SIVALab laboratory between Renault and Heudiasyc. The authors would like to thank Rémy Huet and Vincent Brebion for their technical support in this work.

#### REFERENCES

- B. Missaoui, M. Noizet, and P. Xu, "Map-aided annotation for pole base detection," in *IEEE Intelligent Vehicles Symposium Workshop*, June 2023.
- [2] L. Li, M. Yang, L. Weng, and C. Wang, "Robust localization for intelligent vehicles based on pole-like features using the point cloud," *IEEE Transactions on Automation Science and Engineering*, pp. 1–14, 2021.
- [3] M. Sefati, M. Daum, B. Sondermann, K. D. Kreiskother, and A. Kampker, "Improving vehicle localization using semantic and pole-like landmarks," in *IEEE Intelligent Vehicles Symposium*, Los Angeles, CA, USA, June 2017, pp. 13–19.
- [4] R. Spangenberg, D. Goehring, and R. Rojas, "Pole-based localization for autonomous vehicles in urban scenarios," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, Daejeon, South Korea, Oct. 2016, pp. 2161–2166.
- [5] M. Gouda, A. Shalkamy, X. Li, and K. El-Basyouny, "Fully automated algorithm for light pole detection and mapping in rural highway environment using mobile light detection and ranging point clouds," *Transportation Research Record: Journal of the Transportation Research Board*, vol. 2676, no. 7, pp. 617–629, July 2022.
- [6] M. Lehtomäki, A. Jaakkola, J. Hyyppä, A. Kukko, and H. Kaartinen, "Detection of vertical pole-like objects in a road environment using vehicle-based laser scanning data," *Remote Sensing*, vol. 2, no. 3, pp. 641–664, Feb. 2010.
- [7] B. Rodríguez-Cuenca, S. García-Cortés, C. Ordóñez, and M. Alonso, "Automatic detection and classification of pole-like objects in urban point cloud data using an anomaly detection algorithm," *Remote Sensing*, vol. 7, no. 10, pp. 12680–12703, Sept. 2015.
- [8] F. Ghallabi, G. El-Haj-Shhade, M.-A. Mittet, and F. Nashashibi, "Lidarbased road signs detection for vehicle localization in an HD map," in *IEEE Intelligent Vehicles Symposium*, Paris, France, June 2019, pp. 1484–1490.
- [9] 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 *Computer Vision – ECCV*, D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, Eds. Cham: Springer International Publishing, 2014, pp. 740–755.
- [10] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? The KITTI vision benchmark suite," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2012, pp. 3354–3361.
- [11] J. Behley, M. Garbade, A. Milioto, J. Quenzel, S. Behnke, C. Stachniss, and J. Gall, "SemanticKITTI: A dataset for semantic scene understanding of lidar sequences," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, October 2019.
- [12] 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*, June 2016, pp. 3212–3223.
- [13] F. Yu, H. Chen, X. Wang, W. Xian, Y. Chen, F. Liu, V. Madhavan, and T. Darrell, "BDD100K: A diverse driving dataset for heterogeneous multitask learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, June 2020, pp. 2636–2645.

![](_page_8_Picture_0.jpeg)

(a) Nominal case

(b) Good ground fitting

l fitting (c) U

(c) Unfiltered poles (d) Non-annotated poles

(e) Unmapped features

Fig. 8: Examples of map-aided annotated images from our dataset in different situations. The first row corresponds to a naive projection of the HD map features with a simple filtering of far away features. The second row adds a lidar-based filtering to remove the occluded features. The last row adds an additional ground segmentation to refine the height estimation.

![](_page_8_Picture_7.jpeg)

(a) Ground truth

(b) Prediction from BDD100K model

(c) Prediction from Compiègne model

Fig. 9: Ground truth and predictions from models trained on BDD100K and Compiègne on a Compiègne validation image. The green crosses are true positives, red ones false negatives and the blue circles are misdetections.

- [14] H. Dong, X. Chen, S. Särkkä, and C. Stachniss, "Online pole segmentation on range images for long-term lidar localization in urban environments," *Robotics and Autonomous Systems*, vol. 159, p. 104283, 2023.
- [15] C. Sun, J. M. U. Vianney, Y. Li, L. Chen, L. Li, F.-Y. Wang, A. Khajepour, and D. Cao, "Proximity based automatic data annotation for autonomous driving," *IEEE/CAA Journal of Automatica Sinica*, vol. 7, no. 2, pp. 395–404, 2020.
- [16] W. H. Lee, K. Jung, C. Kang, and H. S. Chang, "Semi-automatic framework for traffic landmark annotation," *IEEE Open Journal of Intelligent Transportation Systems*, vol. 2, pp. 1–12, 2021.
- [17] S. Lee, H. Lim, and H. Myung, "Patchwork++: Fast and robust ground segmentation solving partial under-segmentation using 3d point cloud," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2022, pp. 13 276–13 283.
- [18] C.-Y. Wang, A. Bochkovskiy, and H.-Y. M. Liao, "YOLOv7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors," arXiv preprint:2207.02696, 2022.
- [19] L. Vilalta Estrada, C. Muñoz García, E. Domínguez Tijero, M. Noizet, P. Xu, S. Y. Voon, S. Guerassimov, and W. W. Cox, "ERASMO – Enhanced Receiver for Autonomous MObility," in *Proceedings of the* 15th ITS European Congress, May 2023.