# Mix-Teaching: A Simple, Unified and Effective Semi-Supervised Learning Framework for Monocular 3D Object Detection

Lei Yang[1], Xinyu Zhang[1],[*] Li Wang[1], Minghan Zhu[2], Chuang Zhang[1], Jun Li[1]

[1]State Key Laboratory of Automotive Safety and Energy, Tsinghua University
[2]University of Michigan

{yanglei20, zhch20}@mails.tsinghua.edu.cn; {xyzhang, wangli_thu, lijun1958}@mail.tsinghua.edu.cn
minghanz@umich.edu

## Abstract

*Monocular 3D object detection is an essential perception task for autonomous driving. However, the high reliance on large-scale labeled data make it costly and time-consuming during model optimization. To reduce such over-reliance on human annotations, we propose Mix-Teaching, an effective semi-supervised learning framework applicable to employ both labeled and unlabeled images in training stage. Mix-Teaching first generates pseudo-labels for unlabeled images by self-training. The student model is then trained on the mixed images possessing much more intensive and precise labeling by merging instance-level image patches into empty backgrounds or labeled images. This is the first to break the image-level limitation and put high-quality pseudo labels from multi frames into one image for semi-supervised training. Besides, as a result of the misalignment between confidence score and localization quality, it's hard to discriminate high-quality pseudo-labels from noisy predictions using only confidence-based criterion. To that end, we further introduce an uncertainty-based filter to help select reliable pseudo boxes for the above mixing operation. To the best of our knowledge, this is the first unified SSL framework for monocular 3D object detection. Mix-Teaching consistently improves MonoFlex and GUPNet by significant margins under various labeling ratios on KITTI dataset. For example, our method achieves around +6.34% AP@0.7 improvement against the GUPNet baseline on validation set when using only 10% labeled data. Besides, by leveraging full training set and the additional 48K raw images of KITTI, it can further improve the MonoFlex by +4.65% improvement on AP@0.7 for car detection, reaching 18.54% AP@0.7, which ranks the 1st place among all monocular based methods on KITTI test leaderboard. The code and pretrained models will be released at here.*
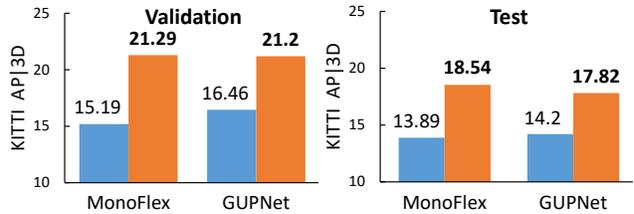
---

[*]Corresponding author: xyzhang@tsinghua.edu.cn



Figure 1. **Performance comparison.** The fully-supervised baselines are represented in cyan, and their semi-supervised training results with our Mix-Teaching are displayed in orange. Our proposed method outperforms the baselines by a large margin on both KITTI validation and test set for Car category.

## 1. Introduction

Monocular 3D object detection is the task of predicting the categories and 3D bounding boxes for surrounding objects with a single image. Owing to its distinct advantages and potential applications in autonomous driving and robotics, this task has attracted extensive attention of researchers from both academia and industry. In recent years, many innovative detectors have emerged and achieved increasing accuracy. However, most of these methods are heavily dependent on labeled data. Compared with the human-annotated images that are often expensive and time-consuming, raw images are easier to achieve large-scale collection. Thus, taking full advantage of both labeled and unlabeled data in model training is a promising approach to alleviate the heavy reliance on human annotations.

Semi-supervised learning (SSL) can help effectively improve the performance of fully-supervised baselines by employing both labeled and unlabeled data. In recent years, plentiful SSL methods for classification [1, 41, 51, 55], 2D object detection [9, 12, 13, 17, 42, 44, 50, 52, 53, 56] and LiDAR-based 3D object detection [46, 59] have been proposed and applied. It can be generally divided into pseudo labeling and consistency regularization. Pseudo labeling

first generates pseudo labels for unlabeled data by self-training [51] or Mean Teacher [45]. Then the student model is trained to predict the same pseudo labels on the same unlabeled images applied with label-preserving data augmentations. In this way, the student model can learn useful information from pseudo labels. Consistency regularization adds a consistency loss to enforce the model make stable predictions on different disturbed data, which helps improving model generalization ability. But as far as we know, there hardly exists a semi-supervised learning framework specially designed for monocular 3D object detection.

Due to the challenge of recovering depth information from a single image, the accuracy of monocular 3D object detection is lagging significantly behind that of 2D object detection and LiDAR-based 3D object detection. Take KITTI [10] benchmark for example, the state of the art methods for monocular 3D object detection just achieve less than 15% AP@0.7, while the candidates for the other two tasks have reached more than 85%-96% AP@0.7. This means that the pseudo labels for unlabeled images are predicted by image-based 3D object detectors with lower precision and recall. Lower precision signifies that more incorrect predictions are possible to be used as labels for unlabeled samples, which can lead to serious confirmation bias. On the other hand, lower recall explains the true positive labels on each image are far from enough to provide adequate supervision signals. Meanwhile, lacking of pseudo labels for plentiful objects can further cause miss-detection. However, most existing SSL methods directly employ the original unlabeled images or the mixup of two images as in Instant-Teaching [9], which can't handle the situations imposed by low recall pseudo labels effectively. To overcome these issues, we propose Mix-Teaching, an general semi-supervised learning framework for most monocular 3D object detectors.

One key challenge for semi-supervised monocular 3D object detection is the extremely low recall of pseudo-labels. For an unlabeled image, only a small part of the objects on this image are correctly detected (less true positive), whereas most of the remaining instances are ignored (more false nagetive). Sparsely distributed true positives fail to provide adequate supervision signals in semi-supervised training. Meanwhile, training data with overmuch missing labels tends to bring miss-detection to most monocular 3D object detectors. In the proposed Mix-Teaching, we firstly predict pseudo-labels for unlabeled data by self-training. Unlabeled samples are then split into image patches collection with high-quality pseudo labels and the collection of background images containing no objects. Subsequently, the student model is trained on the mixed images that are created by merging the above instance image patches into empty backgrounds or human labeled images through strong data augmentation. In this way, the generated im-

ages are full of instances with high-quality pseudo labels while successfully avoiding the missing label cases, which is more effective for semi-supervised training. Finally, we adopt multi-stage training scheme to progressively propagate information from the labeled to the unlabeled data.

Another key challenge for semi-supervised monocular 3D object detection is the confirmation bias. In other words, the model is overfitting to incorrect pseudo labels, which is caused by the extremely low precision of image-based 3D object detectors. Considering the misalignment between confidence score and localization quality, it's not exhaustive to eliminate incorrect labels using only confidence-based filter. To this end, we further propose an uncertainty-based filter to help remove noisy pseudo labels. In this method, the predictions by the models with identical structure but different parameters are used to estimate the uncertainty to each object. For the prediction set belonging to the same object, the higher uncertainty, the less predictions in this set, and the larger localization misalignment among them. We build a formula to represent the localization uncertainty in 3D object detection task. Base on both confidence-based filter and uncertainty-based filters, we manage to remove incorrect pseudo labels more effectively in semi-supervised training and thus alleviate the confirmation bias. Because the process of removing noisy pseudo labels is only carried out at the beginning of each training stage, the efficiency influence from uncertainty calculation is inappreciable.

We benchmark Mix-teaching with SSL setting using the full KITTI [10] object data and KITTI [10] raw data. When using MonoFlex [58] as backbone detector, Mix-Teaching achieve state of the art results on KITTI test leaderboard, which even surpasses the LPCG [33] method that directly relies on LiDAR-based 3D object detectors to generate pseudo labels. Furthermore, we provide the SSL experiments under different labling ratio, which can serve an initial baseline for semi-supervised monocular 3D object detection.

Our contributions can be summarized as follows:

- We clarify the main difficulties in accomplishing semi-supervised learning for monocular 3D object detection and explain why existing SSL approaches can't handle these issues. On this basis, we propose Mix-Teaching, a general semi-supervised framework for monocular 3D object detection.

- To alleviate the confirmation bias, we further propose an uncertainty-based filter to help remove noisy pseudo labels effectively.

- Extensive experiments on KITTI dataset demonstrate the significant efficacy of Mix-Teaching framework. As the first study of SSL for monocular 3D object detection, this can serve as a crucial baseline for further researches.

## 2. Related Work

**Monocular 3D Object Detection.** A number of methods have been proposed for monocular 3D object detection. How to reconstruct spatial information more effectively is the core problem of these approaches. The Pseudo-LiDAR-based methods [8, 28, 29, 49, 54] firstly transform the input image to dense artificial point clouds with the existing depth estimation algorithms [7, 34] and then employs LiDAR-based 3D object detectors [15, 60]. The geometry-based methods [16, 22, 31] infer depth information based on the 2D/3D geometry constraint of specific reference. Another keypoint-based works [3, 5, 25, 27, 30, 36, 58, 61] directly estimate the 3D properties of instance relying on the high-dimensional features at keypoint position.

**Semi-supervised Learning.** Semi-supervised Learning focuses on training models with both labeled and unlabeled data, which has achieved state-of-the-art performance on classification [1, 41, 51, 55], 2D object detection [9, 12, 13, 17, 42, 44, 50, 52, 53, 56] and LiDAR-based 3D object detection [46, 59]. One popular type of SSL is consistency regularization, which constrains the outputs of different augmented inputs to be consistent. CSD [12] is a consistency-based method for 2D object detection. This approach ensures the consistent predictions between input images and their flipped versions on both labeled and unlabeled data. SESS [59] is a semi-supervised learning framework for LiDAR-based 3D object detection. To enhance the model generalization ability, this method applies three consistency losses on two sets of 3D proposals from teacher and student networks. The other kind of SSL is pseudo labeling, which is based on high-quality pseudo labels and can be seen as the hard version of consistency regularization. FixMatch [41] first generates pseudo labels on weakly augmented unlabeled images, and then the student model is trained to predict the same classifications on strong augmented data. Unbiased Teacher [24] addresses the pseudo-labeling bias issue caused by class imbalance in 2D annotations with the help of EMA [45] training and focal loss [21]. 3DIoUMatch [46] achieves semi-supervised 3D object detection in point cloud with a teacher-student mutual learning framework. To improve the quality of pseudo labels, all the predictions that fail to pass the thresholds on classification score, objectness confidence and 3D IoU will be filtered out.

In spite of the success of SSL in classification, 2D object detection and LiDAR-based 3D object detection, there hardly exists a general semi-supervised learning framework specialized for monocular 3D object detection.

## 3. Method

In this section, we first give a mathematical definition of semi-supervised monocular 3D object detection task (see Section 3.1). Then, we show an overview of our training schema (see Section 3.2) and Mix-Teaching framework (see Section 3.3). The uncertainty-based filter is introduced in Section 3.4.

### 3.1. Problem Definition

In semi-supervised monocular 3D object detection, we have labeled data $I^L = \left\{ (x_1^L, y_1^L), \ldots, (x_{n_l}^L, y_{n_l}^L) \right\}$ and abundant unlabeled data $I^U = \left\{ x_1^U, \ldots, x_{n_u}^U \right\}$, where $x$ is image, $y$ denotes the human-annotated label that contains category and 3D bounding box. $n_l$ and $n_u$ represent the number of labeled and unlabeled images respectively. We aim to significantly improve the performance of fully-supervised baselines by applying both labeled and unlabeled data in training.

### 3.2. Multi-stage Training Schema

We adopt a multi-stage training schema. The initial teacher model is trained on labeled data, followed by a pseudo-labeling process for the unlabeled data. Then we train a noisy student model using all the labeled and unlabeled images following decomposition and re-combination technique. This resulting model will be used as a new teacher model in the next stage.

### 3.3. Mix-Teaching Framework

We propose a SSL framework for monocular 3D object detection, called Mix-Teaching, as shown in Figure 2. This is a general approach that can be easily applied to most monocular 3D object detectors. Our Mix-Teaching is mainly composed of two stages: database-oriented pseudo-labeling and noisy student with mixed data.

**Database-oriented Pseudo Labeling.** To make the most of sparsely distributed high-quality pseudo labels in semi-supervised training, all the labels and background images need to be gathered together. As shown in Figure 2, we perform a test-time inference of the teacher model on unlabeled images to generate pseudo labels. By applying confidence-based and uncertainty-based filters, we create an instance database that is composed of instance-level image patches and their corresponding high-quality pseudo labels. Based on the object existence filter, we select all the background images that don't contain any predictions from unlabeled data and create a background database

**Noisy Student with Mixed Data.** Based on the above two databases and labeled images, we create the mixed images containing more intensive and precise labels for semi-supervised training. There are two general strategies for this purpose. One way is to paste the image patches from the instance database on labeled images. Another way is to paste the instance-level patches on the images that come from the background database. During the process, the instance-level patches are pasted to target images according to their 2D bounding box on source images. To avoid over occlusion
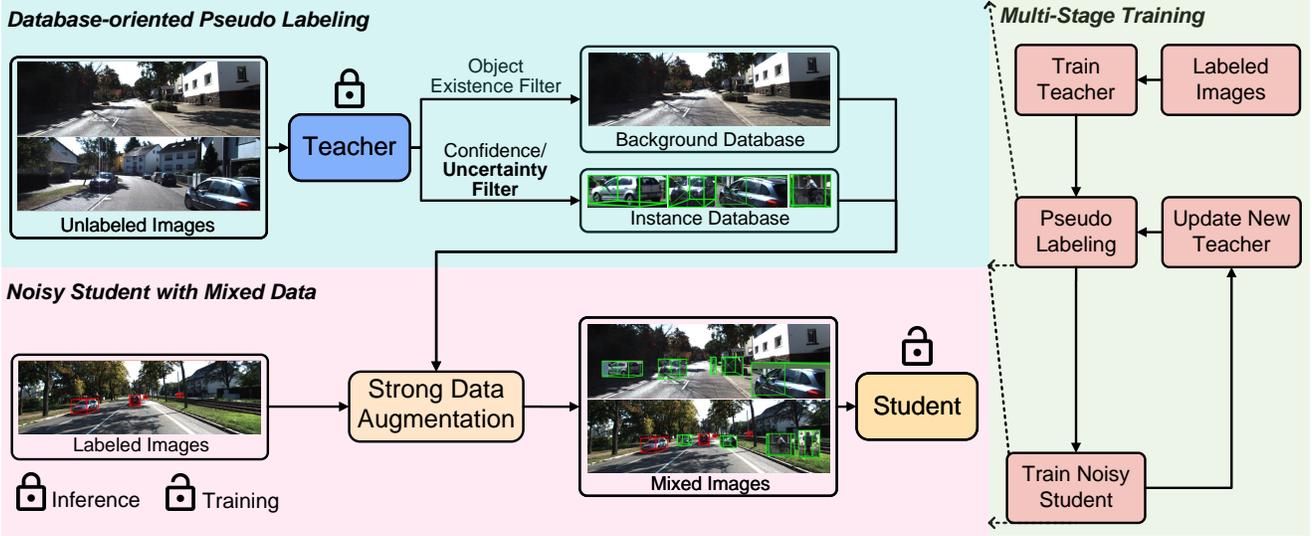
Figure 2. **Overview of Mix-Teaching Framework.** Mix-Teaching follows multi-stage training scheme. There are two crucial process in each training stage. **Database-oriented Pseudo Labeling:** Based on the pseudo labels that are generated by applying the teacher model on unlabeled images, We create two databases: one is background database that consists of images without any objects, the other one is instance database composed of image patches with high-quality pseudo labels. **Noisy Student with Mixed Data:** The student model is trained on the mixed images that possess much more intensive and precise labels by merging image patches from the above instance database to labeled images or the ones from background database. See Section 3.2 and Section 3.3 for more details.

and other impossible outcomes, we additionally perform a 2D bounding box collision test to remove invalid paste operations.

To alleviate the confirmation bias and improve the model generalization ability, we further propose a series of box-level strong data augmentations. For completeness, we describe the list of augmentations below. Each operation has a magnitude that decides the augmentation degree of strength. We visualize the augmented instances with single or fusion strategies mentioned above in Figure 3.

1. Border Cut (**B**): Before the pasting operation, cutting the horizontal or vertical border of image patch with a random ratio (0-0.3).

2. Color Padding(**C**): Similar to the border cut, but replacing the cut operation with random color padding.

3. Mixup(**M**): making a weighted average between the foreground image patches and backgrounds with a random ratio (0.6-1.0).

When given a batch of mixed images and the corresponding human labels or pseudo labels, the model is trained by jointly minimizing the supervised loss and unsupervised loss as follows:

$$\mathcal{L} = \mathcal{L}_s + \lambda \times \mathcal{L}_u \tag{1}$$

where $\mathcal{L}$ is total loss, we use hyper-parameter $\lambda$ to balance the supervised loss $\mathcal{L}_s$ and the unsupervised loss $\mathcal{L}_u$.

The supervised loss $\mathcal{L}_s$ consists of a classification loss $\mathcal{L}_{cls}$ and a regression loss $\mathcal{L}_{reg}$. It can be calculated as:

$$\mathcal{L}_s = \sum_L \frac{1}{N_l} \sum_i (\mathcal{L}_{cls}(b_l^i) + \mathcal{L}_{reg}(b_l^i)) \tag{2}$$

where $L$ denotes the index of labeled images in a batch, $N_l$ represents the number of human annotations for each image, $b_l^i$ is the $i$-th label in the $L$-th labeled image.

The unsupervised loss $\mathcal{L}_u$ is computed on pseudo labels and can be written as:

$$\mathcal{L}_u = \sum_L \frac{1}{N_u} \sum_i (\mathcal{L}_{cls}(b_u^i) + \mathcal{L}_{reg}(b_u^i))$$
$$+ \sum_B \frac{1}{N_u} \sum_i (\mathcal{L}_{cls}(b_u^i) + \mathcal{L}_{reg}(b_u^i)) \tag{3}$$

where $B$ indicates the index of background images, $N_u$ is the number of pseudo labels on each image, $b_u^i$ represents the $i$-th pseudo label on a labeled or a background image.

For monocular 3D object detection, the extremely low recall and precision of related detectors make it a great challenge to apply the existing semi-supervised methods for 2D object detection that focus more on false positives but ignore false negatives to this field. In the Mix-Teaching proposed above, following decomposition and recombination methodology, we collect all positive instances and merge them into backgrounds to create newly mixed images for

Figure 3. **Visualization of box-level strong augmentations.** From left to right: original image patch, border cut, color padding, mixup and the fusion of previous three methods.

semi-supervised training. The positive instances indicate object-level image patches with high-quality pseudo labels. The backgrounds denote empty unlabeled images or labeled data. The newly mixed images will possess high recall and less false positives at the same time, which is effective to solve the extremely low recall and confirmation bias challenges in semi-supervised monocular 3D object detection.

### 3.4. Uncertainty-based Filter

As shown in Figure 4(a), there exists a huge misalignment between the classification score and the localization precision of box candidates. A considerable proportion of predictions have a high confidence score but low 3D IoU with ground truth. When the high-quality pseudo labels are discriminate completely based on the confidence-based filter, a lot of incorrect pseudo labels will be used for semi-supervised training, which will strengthen the confirmation bias.

In order to alleviate the above issues, it is necessary to remove noisy labels in semi-supervised training. To this end, we further propose an uncertainty-based filter in which we infer localization uncertainty on the basis of the discrepancy of the predictions for one object from multi models.

When a certain image is given to $N$ isomorphic models with different parameters. For a specific object on this image, there will be $M$ predictions. We define the localization uncertainty mainly from two points of view: (1) the number of predictions $M$ associated with this object; (2) The discrepancy between these predicted boxes. The predictions number $M$ reflects the level of missed detection among $N$ models. The disparity in box candidates reveals the randomness in model predictions.

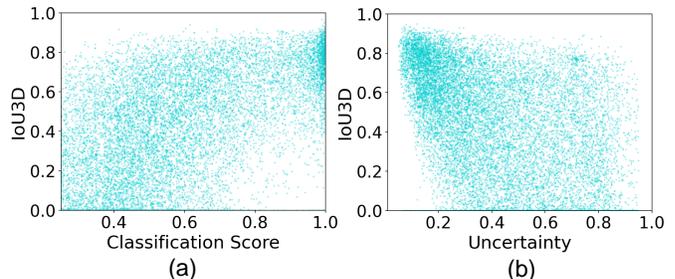We calculate the uncertainty in the following steps:



Figure 4. **The statistics of MonoFlex predictions on KITTI validation set.** (a) the relationship between the 3D IoU with ground truth box and classification score. (b) the relationship between the 3D IoU and localization uncertainty.

1. All predictions from N models are stored in list $B$.

2. Declare three lists $G$, $H$ and $U$. $G$ is used to store box clusters. Each cluster represents the predictions for a certain object from $N$ models . $H$ is for the box with highest confidence score in each cluster. $U$ saves the localization uncertainty for each box in list $H$.

3. Iterate through all the boxes in list $B$ to find the matching box that belongs to the current cluster $C$. The matching condition is defined as a box with a large overlap with the initial box $b_m$ of cluster $C$ under the condition $IoU3D > thr$. All matching boxes will be moved from list $B$ to cluster $C$. And then update the current cluster $C$ to list $G$.

4. If there are still unprocessed boxes in list B, select the box $b_m$ that has the maximum score in list $B$ and move it to list $H$. Initialize a new cluster $C$ with box $b_m$ and

proceed to step 3.

5. When all boxes in B are processed, calculate the uncertainty $u$ for each box cluster $C$ in list $G$ with the following equations. The results are added to list $U$.

$$u = uncertain(C) =$$

$$1 - \frac{\sum_{i=0}^{M} \sum_{j=0}^{M} a_{ij} \times IoU3D(b_i, b_j)}{\sum_{i=0}^{N} \sum_{j=0}^{N} a_{ij}} \quad (4)$$

$$a_{ij} = \begin{cases} 1 & if \ i \neq j \\ \beta & if \ i = j \end{cases}, \quad (5)$$

where $M$ is the number of boxes in cluster $C$, $N$ denotes the number of models, $b_i$ is the $i$-th box in cluster $C$, $a_{ij}$ represents the weight for each item. $\beta$ is the hyper-parameter that controls how the number of box candidates affects uncertainty.

The uncertainty $u$ ranges from 0 to 1. When the value is 0, it indicates that there exists no miss-detection in $N$ models ($M = N$), and all $N$ box candidates are perfectly consistent. When the value is 1, it means that all models fail to detect this object.

As shown in Figure 4 (b), we visualize the relationship between the 3D IoU and localization uncertainty. Compared with the classification score, the obtained uncertainty can better measure the localization accuracy.

Compared with 3D confidence [39] or 3D IoU [18, 46] that requires related branch design for specific detectors, our uncertainty-based filter is model-independent and can be applied to many types of image-based 3D object detectors, which is more appropriate for the proposed general semi-supervised learning framework.

## 4. Experiments

### 4.1. Dataset and Metrics

Contrasted with nuscenes [2] and waymo [43] dataset that lacks of a large amount of unlabeled data, KITTI [10] dataset provides 15K frames labeled data and 48K unlabeled images, which is more appropriate for semi-supervised learning research that relies on limited labeled data and larger scale unlabeled images. Therefore, we evaluate our Mix-teaching on the challenging KITTI [10] dataset. KITTI contains 7,481 images for training and 7,518 images for testing. Since we have no access to the manual annotations of testing set, the training set is further split into 3,712 training samples and 3,769 validation samples as mentioned in [4] for local evaluation. Besides, there are additional raw data consists of 48K temporal images. These images don't coincide with the above training or testing set, and thus can be used as unlabeled data for semi-supervised

training. We use the average precision on Car, Pedestrian and Cyclist for 3D and bird's eye view (BEV) object detection as the metrics. Following [40], all the evaluation results on validation and testing set are based on $AP_{40}$ instead of the original 11-point interpolated average precision.

### 4.2. Implement Details

We adopt the MonoFlex [58] and GUPNet [27] as two baseline detectors. The localization uncertainty is calculated based on five models from different training rounds. During the pseudo labeling process, only the predictions with confidence score larger than 0.7 and localization uncertainty less than 0.25 will be added to the instance database. The images without any detections are collected to build background database. During the student model training period, We initialize the student with the previous teacher model. the images from background database are selected with a chance of 42% apart from labeled images. For the box-level strong data augmentation, we apply mixup on every instance image patches. border cut and color padding augmentation are employed with a chance of 50%. We set the hyper-parameter $\lambda = 1.0$. Following the multi-training scheme, we conduct three cycles of semi-supervised training for all experiments.

### 4.3. Quantitative Results

**Comparison with Fully-supervised Baselines.** We make a detailed comparison with the supervised baselines, including GUPNet [27] and MonoFlex [58], under different ratios of training set. All 48K raw images of KITTI are used as unlabeled data for semi-supervised training. As depicted in Table 1, Mix-Teaching significantly outperforms MonoFlex [58] and GUPNet [27] under each ratio settings, which verify the effectiveness of our semi-supervised framework. when using only 10% labeled data, our approach gains around +6.34% and +5.98% $AP_{3D}$ improvements on moderate level over MonoFlex and GUPNet baselines. This indicates our framework is able to learn knowledge from unlabeled data, and the effect is more obvious when the number of labeled data is scarce. Furthermore, it is worth pointing out that when using all training set, our Mix-Teaching is able to further outperforms the upper-bound performance of two baselines by a large margin.

**Results on KITTI Test Set.** We evaluate the proposed Mix-Teaching on KITT test set using MonoFlex [58] and GUPNet [27] as two base monocular detectors. Table 2 shows the quantitative results of our method and other top performance detectors from the official KITTI leaderboard. Overall, Mix-Teaching achieves superior results over all two baselines across all settings under fair conditions. For instance, the proposed method improve the $AP_{3D}$ of GUPNet [27] by **+7.44/+3.62/+2.95** absolute improvements under easy/moderate/hard setting. Meanwhile,

| Method | 10% | | | 50% | | | 100% | | |
| | Easy | Mod. | Hard | Easy | Mod. | Hard | Easy | Mod. | Hard |
|---|---|---|---|---|---|---|---|---|---|
| GUPNet | 8.42 / 13.61 | 5.14 / 8.67 | 4.11 / 7.35 | 18.21 / 26.28 | 14.14 / 19.38 | 11.12 / 16.63 | 22.76 / 31.07 | 16.46 / 22.94 | 13.72 / 19.75 |
| Ours | 16.55 / 22.57 | 11.48 / 15.90 | 9.44 / 13.13 | 25.67 / 34.87 | 19.76 / 26.19 | 17.07 / 21.76 | 29.12 / 38.48 | 21.04 / 28.25 | 17.56 / 24.36 |
| Abs. Imp. | +8.13 / 8.96 | +6.34 / +7.23 | +5.33 / +5.78 | +7.46 / +8.59 | +5.62 / +6.91 | +4.95 / +5.13 | +6.36 / +7.41 | +4.58 / +5.31 | +3.84 / +4.61 |
| MonoFlex | 5.76 / 9.93 | 4.67 / 7.68 | 3.54 / 6.09 | 21.91 / 28.24 | 15.43 / 20.51 | 13.09 / 18.37 | 23.64 / 29.86 | 17.51 / 23.05 | 14.83 / 20.68 |
| Ours | 14.43 / 18.80 | 10.65 / 14.28 | 8.41 / 11.73 | 29.34 / 36.23 | 20.63 / 26.70 | 17.31 / 23.07 | 29.74 / 37.45 | 22.27 / 28.99 | 19.04 / 25.31 |
| Abs. Imp. | +8.57 / +8.87 | +5.98 / +6.60 | +4.87 / +5.64 | +7.43 / +7.99 | +5.20 / +6.19 | +4.22 / +4.70 | +6.10 / +7.56 | +4.76 / +5.94 | +4.21 / +4.63 |

Table 1. **Quantitative results of $AP_{3D}/AP_{BEV}(IoU=0.7)|_{R_{40}}$ on KITTI val set under under different ratios of training set.** "Abs. Imp." represents absolute improvements.

| Method | Reference | GPU | Runtime (ms) | $AP_{3D}(IoU=0.7)|_{R_{40}}$ | | | $AP_{BEV}(IoU=0.7)|_{R_{40}}$ | | |
| | | | | Easy | Mod. | Hard | Easy | Mod. | Hard |
|---|---|---|---|---|---|---|---|---|---|
| MonoGRNet [35] | TPAMI 2021 | Tesla P40 | 60 | 9.61 | 5.74 | 4.25 | 18.19 | 11.17 | 8.73 |
| MonoPair [5] | CVPR 2020 | Tesla V100 | 60 | 13.04 | 9.99 | 8.65 | 19.28 | 14.83 | 12.89 |
| RTM3D [19] | ECCV 2020 | 1080Ti | 50 | 14.41 | 10.34 | 8.77 | 19.17 | 14.20 | 11.99 |
| D⁴LCN [8] | CVPR 2020 | 1080Ti | 200 | 16.65 | 11.72 | 9.51 | 22.51 | 16.02 | 12.55 |
| Monodle [30] | CVPR 2021 | 1080Ti | 40 | 17.23 | 12.26 | 10.29 | 24.79 | 18.89 | 16.00 |
| MonoRUn [3] | CVPR 2021 | 1080Ti | 70 | 19.65 | 12.30 | 10.58 | 27.94 | 17.34 | 15.24 |
| GrooMeD-NMS [14] | CVPR 2021 | Titian X | 120 | 18.10 | 12.32 | 9.65 | 26.19 | 18.27 | 14.05 |
| MonoRCNN [38] | ICCV 2021 | - | 70 | 18.36 | 12.65 | 10.03 | 25.48 | 18.11 | 14.10 |
| DDMP-3D [47]∗ | CVPR 2021 | - | 180 | 19.71 | 12.78 | 9.80 | 28.08 | 17.89 | 13.44 |
| Ground-Aware [23] | RAL 2021 | 1080Ti | 50 | 21.65 | 13.25 | 9/91 | 29.81 | 17.98 | 13.08 |
| PCT [48] | NIPS 2021 | - | 487 | 21.00 | 13.37 | 11.31 | 29.65 | 19.03 | 15.92 |
| CaDDN [37] | CVPR2021 | 2080Ti | 485 | 19.17 | 13.41 | 11.46 | 27.94 | 18.91 | 17.19 |
| DFR-Net [63]∗ | ICCV 2021 | - | 180 | 19.40 | 13.63 | 10.35 | 28.17 | 19.17 | 14.84 |
| MonoEF [62] | TPAMI 2021 | - | 30 | 21.29 | 13.87 | 11.71 | 29.03 | 19.70 | 17.26 |
| MonoFlex [58]† | CVPR 2021 | 2080Ti | 30 | 19.94 | 13.89 | 12.07 | 28.23 | 19.75 | 16.89 |
| AutoShape [26] | ICCV 2021 | 2080Ti | 52 | 22.47 | 14.17 | 11.36 | 30.66 | 20.08 | 15.95 |
| GUPNet [27]† | ICCV 2021 | 2080Ti | 26 | 20.11 | 14.20 | 11.77 | - | - | - |
| MonoDTR [11] | CVPR 2022 | - | 37 | 21.99 | 15.39 | 12.73 | 28.59 | 20.38 | 17.14 |
| MonoDETR [57] | CVPR 2022 | - | 40 | **23.65** | 15.92 | 12.99 | 32.08 | 21.44 | 17.85 |
| MonoDistill [6]∗ | ICLR 2022 | 1080 Ti | 40 | 22.97 | 16.03 | 13.60 | 31.87 | 22.59 | 19.72 |
| MonoJSG [20] | CVPR 2022 | - | 42 | 24.69 | 16.14 | 13.64 | 32.59 | 21.26 | 18.18 |
| DD3D [32]∗ | ICCV 2021 | 2080Ti | 60 | 23.19 | 16.87 | 14.36 | 32.35 | 23.41 | 20.42 |
| LPCG [33]∗ | Arxiv 2021 | 2080Ti | 30 | 25.56 | 17.80 | <u>15.38</u> | <u>35.96</u> | **24.81** | **21.86** |
| **GUPNet + Ours** | - | 2080Ti | 26 | **27.55** | <u>17.82</u> | 14.72 | **36.39** | 24.14 | <u>20.49</u> |
| Abs. Imp. | - | - | - | +7.44 | +3.62 | +2.95 | - | - | - |
| **MonoFlex + Ours** | - | 2080Ti | 30 | <u>26.89</u> | **18.54** | **15.79** | 35.74 | <u>24.23</u> | <u>20.80</u> |
| Abs. Imp. | - | - | - | +6.95 | +4.65 | +3.72 | +7.51 | +4.48 | +3.91 |

Table 2. **Performance of the Car category on KITTI test set**. We use **bold** to highlight the highest results and <u>underlined</u> for the second-highest ones. † represents the baseline we employed. All methods are ranked by $AP_{3D}$ on moderate setting (same as KITTI leaderboard), Our method outperforms the baseline by a large margin and achieves the best performance.

our approach increases the same metric of MonoFLex [58] from 19.94/13.89/12.07 to 26.89/18.54/15.79, which is absolutely remarkable. What's more, Our Mix-Teaching even surpasses the LPCG [33] that directly relies on LiDAR-based 3D object detectors to help generate pseudo labels on $AP_{3D}$ metric using the same MonoFlex baseline. We rank the 1st place according to $AP_{3D}$ on moderate setting (same as KITTI leaderboard).

**Results on Pedestrian and Cyclist Categories.** Compared with Car, Pedestrian and Cyclist are more challenging to be detected owing to their non-rigid structure, small scale. As shown in Table 3, our Mix-teaching can further boost the $AP_{3D}(IoU=0.5)|_{R_{40}}$ metric of MonoFlex [58] baseline around 18% relative improvements for pedestrian and 108% for cyclist on the test set, which demonstrates the effectiveness of our method on small-scale objects.

| Method | Cat. | $\text{AP}_{3D}\|_{R_{40}}$ | | | $\text{AP}_{BEV}\|_{R_{40}}$ | | |
|---|---|---|---|---|---|---|---|
| | | Easy | Mod. | Hard | Easy | Mod. | Hard |
| MonoFlex | | 9.43 | 6.31 | 5.26 | 10.36 | 7.36 | 6.29 |
| Ours | Ped. | **11.67** | **7.47** | **6.61** | **12.34** | **8.40** | **7.06** |
| Rel. Imp.(%) | | 23.75 ↑ | 18.38 ↑ | 25.67 ↑ | 19.11 ↑ | 14.13 ↑ | 12.24 ↑ |
| MonoFlex | | 4.17 | 2.35 | 2.04 | 4.41 | 2.67 | 2.50 |
| Ours | Cycl. | **8.04** | **4.91** | **4.15** | **8.56** | **5.36** | **4.62** |
| Rel. Imp.(%) | | 92.81 ↑ | 108.94 ↑ | 103.43 ↑ | 94.10 ↑ | 100.75 ↑ | 84.80 ↑ |

Table 3. **Quantitative results for Pedestrian and Cyclist on KITTI test set.** "Rel. Imp." represents relative improvements.

## 4.4. Ablation Studies

In this section, we perform ablation studies to investigate the effects of each elements. We use MonoFlex [58] as the base detector. The training of ablation experiments is conducted on the full KITTI training set. The results for car category are evaluated on the corresponding validation set.
**The Scale of Unlabeled Data** We investigate if the proposed semi-supervised learning strategy keeps improving performance with increasing unlabeled data. As show in table 4, Compared with the results of 24K KITTI raw data, the experiment when using the whole 48K unlabeled data can further improve the AP3$D$ for car category from 20.61% to 22.27%. This means that, with more unlabeled images, our Mix-Teaching can improve the accuracy of fully supervised detectors to a new level.
**Background Database** Next, we investigate if the backbone database is necessary during the student model training period. As shown in table 4, when using half KITTI raw data, we gain +2.23% and +2.45% absolute improvement over the without background database version on $\text{AP}_{3D}$ and $\text{AP}_{BEV}$ respectively. And when it comes to all the unlabeled data condition, background database brings about significant improvements as well.

| Raw Data | Background | $\text{AP}_{3D}\|_{R_{40}}$ | | | $\text{AP}_{BEV}\|_{R_{40}}$ | | |
|---|---|---|---|---|---|---|---|
| | | Easy | Mod. | Hard | Easy | Mod. | Hard |
| - | - | 23.64 | 17.51 | 14.83 | 29.86 | 23.05 | 20.68 |
| 50% | - | 23.79 | 18.29 | 15.66 | 32.55 | 24.15 | 21.58 |
| | √ | 27.49 | 20.61 | 17.68 | 36.19 | 26.60 | 23.13 |
| 100% | - | 24.22 | 19.02 | 16.16 | 33.33 | 25.70 | 22.34 |
| | √ | **29.74** | **22.27** | **19.04** | **37.45** | **28.99** | **25.31** |

Table 4. **Ablation study on the effects of background database and unlabeled data scale.**

**Box-level Data Augmentations.** We ablate the effects of box-level data augmentations. As shown in Table 5, all border cut, color padding and mixup are helpful to improve performance. The combination of the above three data augmentations can further improve the performance of Mix-Teaching.

| B | C | M | $\text{AP}_{3D}\|_{R_{40}}$ | | | $\text{AP}_{BEV}\|_{R_{40}}$ | | |
|---|---|---|---|---|---|---|---|---|
| | | | Easy | Mod. | Hard | Easy | Mod. | Hard |
| - | - | - | 26.44 | 20.04 | 17.18 | 35.04 | 25.87 | 22.54 |
| √ | - | - | 28.40 | 21.05 | 18.15 | 36.73 | 26.95 | 23.48 |
| - | √ | - | 27.22 | 20.76 | 17.85 | 35.93 | 26.71 | 23.26 |
| - | - | √ | 27.74 | 21.21 | 18.32 | 36.99 | 27.08 | 23.54 |
| √ | √ | √ | **29.74** | **22.27** | **19.04** | **37.45** | **28.99** | **25.31** |

Table 5. **Ablation study on the effects of box-level data augmentations.** "B" implies the border cut, "C" denotes the color padding, "M" represents the mixup.

**Confidence-based and Uncertainty-based Filters.** As shown in Table 6, the results of applying uncertainty-based filter surpass that of only using confidence-based filter by 1.06% $\text{AP}_{3D}$ on moderate set, which explains that our proposed uncertainty-filter is much more effective. When employing both of the two filters, we achieve the best results.

| Conf. | Unc. | $\text{AP}_{3D}\|_{R_{40}}$ | | | $\text{AP}_{BEV}\|_{R_{40}}$ | | |
|---|---|---|---|---|---|---|---|
| | | Easy | Mod. | Hard | Easy | Mod. | Hard |
| - | - | 23.83 | 17.81 | 15.13 | 30.52 | 23.14 | 19.96 |
| √ | - | 26.51 | 20.08 | 17.21 | 35.14 | 25.95 | 22.59 |
| - | √ | 27.78 | 21.14 | 18.18 | 35.36 | 26.87 | 23.52 |
| √ | √ | **29.74** | **22.27** | **19.04** | **37.45** | **28.99** | **25.31** |

Table 6. **Ablation study on the effects of two filters.** "Conf." denotes the filter based on the confidence. "Unc." represents the uncertainty-based filter.

**The Thresholds for Confidence-based and Uncertainty-based Filters.** We studies the effects of different confidence score thresholds and uncertainty thresholds in discriminating high-quality pseudo labels. As shown in Table 7, The best performance is achieved when the confidence threshold is set to 0.7 and the uncertainty threshold is set to 0.25.

| Conf. Thre. | Unc. Thre. | $\text{AP}_{3D}\|_{R_{40}}$ | | | $\text{AP}_{BEV}\|_{R_{40}}$ | | |
|---|---|---|---|---|---|---|---|
| | | Easy | Mod. | Hard | Easy | Mod. | Hard |
| 0.6 | 0.25 | 27.84 | 21.63 | 18.26 | 36.04 | 26.51 | 23.49 |
| 0.7 | 0.25 | **29.74** | **22.27** | **19.04** | **37.45** | **28.99** | **25.31** |
| 0.8 | 0.25 | 28.07 | 21.96 | 18.54 | 36.69 | 27.06 | 23.90 |
| 0.9 | 0.25 | 26.31 | 20.20 | 16.59 | 33.59 | 24.58 | 21.24 |
| 0.7 | 0.45 | 26.97 | 21.65 | 18.31 | 34.09 | 26.50 | 23.45 |
| 0.7 | 0.35 | 28.89 | 21.98 | 18.51 | 36.75 | 26.85 | 23.72 |
| 0.7 | 0.15 | 27.34 | 21.42 | 18.37 | 35.45 | 26.32 | 23.50 |

Table 7. **Ablation study on the effects of different uncertainty and confidence thresholds.** "Conf. Threshold" denotes the threshold of confidence-based filter. "Unc. Threshold" represents the threshold of uncertainty-based filter.
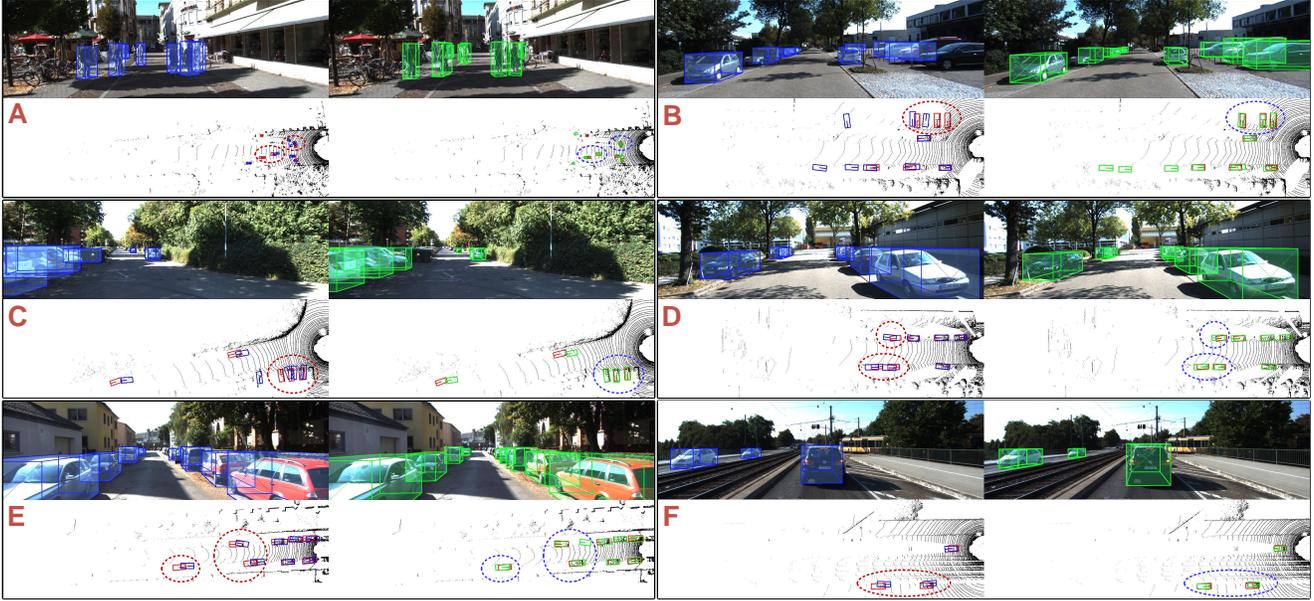
Figure 5. **Qualitative results on the KITTI val set.** We present four pairs of comparisons marked with capital letters from A to F. Each pair consists of four pictures, the upper left displays the predictions of MonoFlex [58] baseline(blue), the lower left is its representation in the bird's-eye view. The upper right shows the results of our Mix-Teaching(green), the lower right is its bird's-eye view display. The red boxes in the bird's eye view represent ground truths. We use dashed ovals to highlight the pronounced difference in the predictions.

## 4.5. Qualitative Analysis

From the qualitative results shown in Figure 5 , Mix-teaching can improve the performance of MonoFlex [58] in various street scenes. As highlighted by the red ovals, our method can produce superior performance for person and cyclist(A in Figure 5), over-occluded objects (C in Figure 5) and cars in both close range (B-C in Figure 5) and long range (D,E,F in Figure 5), which demonstrates the efficiency of our Mix-teaching.

## 5. Conclusion and Future Work

In this paper, we proposed Mix-Teaching, a general semi-supervised learning framework for monocular 3D object detection. Our method first generates pseudo labels for unlabeled data by self-training. Then, following decomposition and re-combination technique, we break the limitation of original images and create newly diverse and label-rich mixed images for semi-supervised training, which can effectively handle the issues imposed by the extremely lower precision and recall of initial pseudo labels. With the proposed uncertainty-based filter, we manage to filter poorly positioned pseudo labels effectively, leading to less noise so as to alleviate confirmation bias. Experiments on KITTI dataset show that Mix-Teaching manages to improve the baseline model by a large margin under various labeling ratios. More importantly, when using 100% training set and MonoFlex as backbone detector, we successfully

rank the first place among all monocular 3D object detectors on KITTI test leaderboard. In this way, we can continuously boost monocular 3D object detectors by collecting more unlabeled images, which has great economic significance in autonomous driving. Moreover, the proposed Mix-Teaching follows the multi-stage training scheme. Adopting end-to-end training fashion will be left for future work.

## References

[1] David Berthelot, Nicholas Carlini, Ekin Dogus Cubuk, Alexey Kurakin, Kihyuk Sohn, Han Zhang, and Colin Raffel. Remixmatch: Semi-supervised learning with distribution matching and augmentation anchoring. In *ICLR*, 2020. 1, 3

[2] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *CVPR*, 2020. 6

[3] Hanshen Chen, Yuyao Huang, Wei Tian, Zhong Gao, and Lu Xiong. Monorun: Monocular 3d object detection by reconstruction and uncertainty propagation. In *CVPR*, pages 10374–10383, 2021. 3, 7

[4] Xiaozhi Chen, Kaustav Kundu, Yukun Zhu, Andrew G. Berneshawi, Huimin Ma, Sanja Fidler, and Raquel Urtasun. 3d object proposals for accurate object class detection. In *NeurIPS*, 2015. 6

[5] Yongjiang Chen, Lei Tai, Kai Sun, and Mingyang Li. Monopair: Monocular 3d object detection using pairwise spatial relationships. In *CVPR*, pages 12090–12099, 2020. 3, 7

[6] Zhiyu Chong, Xinzhu Ma, Hong Zhang, Yuxin Yue, Haojie Li, Zhihui Wang, and Wanli Ouyang. Monodistill: Learning spatial features for monocular 3d object detection. In *ICLR*, 2022. 7

[7] Raúl Díaz and Amit Marathe. Soft labels for ordinal regression. In *CVPR*, pages 4733–4742, 2019. 3

[8] Mingyu Ding, Yuqi Huo, Hongwei Yi, Zhe Wang, Jianping Shi, Zhiwu Lu, and Ping Luo. Learning depth-guided convolutions for monocular 3d object detection. In *CVPRW*, pages 4306–4315, 2020. 3, 7

[9] Qiang feng Zhou, Chaohui Yu, Zhibin Wang, Qi Qian, and Hao Li. Instant-teaching: An end-to-end semi-supervised object detection framework. In *CVPR*, 2021. 1, 2, 3

[10] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, pages 3354–3361, 2012. 2, 6

[11] Kuan-Chih Huang, Tsung-Han Wu, Hung-Ting Su, and Winston H. Hsu. Monodtr: Monocular 3d object detection with depth-aware transformer. In *CVPR*, volume abs/2203.10981, 2022. 7

[12] Jisoo Jeong, Seungeui Lee, Jeesoo Kim, and Nojun Kwak. Consistency-based semi-supervised learning for object detection. In *NeurIPS*, 2019. 1, 3

[13] Jisoo Jeong, Vikas Verma, Minsung Hyun, Juho Kannala, and Nojun Kwak. Interpolation-based semi-supervised learning for object detection. In *CVPR*, 2021. 1, 3

[14] Abhinav Kumar, Garrick Brazil, and Xiaoming Liu. Groomed-nms: Grouped mathematically differentiable nms for monocular 3d object detection. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8969–8979, 2021. 7

[15] Alex H. Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *CVPR*, pages 12689–12697, 2019. 3

[16] Buyu Li, Wanli Ouyang, Lu Sheng, Xingyu Zeng, and Xiaogang Wang. Gs3d: An efficient 3d object detection framework for autonomous driving. In *CVPR*, pages 1019–1028, 2019. 3

[17] Hengduo Li, Zuxuan Wu, Abhinav Shrivastava, and Larry Davis. Rethinking pseudo labels for semi-supervised object detection. *ArXiv*, abs/2106.00168, 2021. 1, 3

[18] Peixuan Li and Huaici Zhao. Monocular 3d detection with geometric constraint embedding and semi-supervised training. *IEEE Robotics and Automation Letters*, 6:5565–5572, 2021. 6

[19] Peixuan Li, Huaici Zhao, Pengfei Liu, and Feidao Cao. Rtm3d: Real-time monocular 3d detection from object keypoints for autonomous driving. In *ECCV*, 2020. 7

[20] Qing Lian, Peiliang Li, and Xiaozhi Chen. Monojsg: Joint semantic and geometric cost volume for monocular 3d object detection. In *CVPR*, 2022. 7

[21] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, pages 2980–2988, 2017. 3

[22] Lijie Liu, Jiwen Lu, Chunjing Xu, Qi Tian, and Jie Zhou. Deep fitting degree scoring network for monocular 3d object detection. In *CVPR*, pages 1057–1066, 2019. 3

[23] Yuxuan Liu, Yuan Yixuan, and Ming Liu. Ground-aware monocular 3d object detection for autonomous driving. *IEEE Robotics and Automation Letters*, pages 919–926, 2021. 7

[24] Yen-Cheng Liu, Chih-Yao Ma, Zijian He, Chia-Wen Kuo, Kan Chen, Peizhao Zhang, Bichen Wu, Zsolt Kira, and Péter Vajda. Unbiased teacher for semi-supervised object detection. In *ICLR*, 2021. 3

[25] Zechen Liu, Zizhang Wu, and Roland T'oth. Smoke: Single-stage monocular 3d object detection via keypoint estimation. In *CVPRW*, pages 4289–4298, 2020. 3

[26] Zongdai Liu, Dingfu Zhou, Feixiang Lu, Jin Fang, and Liangjun Zhang. Autoshape: Real-time shape-aware monocular 3d object detection. In *ICCV*, 2021. 7

[27] Yan Lu, Xinzhu Ma, Lei Yang, Tianzhu Zhang, Yating Liu, Q. Chu, Junjie Yan, and Wanli Ouyang. Geometry uncertainty projection network for monocular 3d object detection. In *ICCV*, 2021. 3, 6, 7

[28] Xinzhu Ma, Shinan Liu, Zhiyi Xia, Hongwen Zhang, Xingyu Zeng, and Wanli Ouyang. Rethinking pseudo-lidar representation. In *ECCV*, 2020. 3

[29] Xinzhu Ma, Zhihui Wang, Haojie Li, Wanli Ouyang, and Pengbo Zhang. Accurate monocular 3d object detection via color-embedded 3d reconstruction for autonomous driving. In *ICCV*, pages 6850–6859, 2019. 3

[30] Xinzhu Ma, Yinmin Zhang, Dan Xu, Dongzhan Zhou, Shuai Yi, Haojie Li, and Wanli Ouyang. Delving into localization errors for monocular 3d object detection. In *CVPR*, pages 4719–4728, 2021. 3, 7

[31] Arsalan Mousavian, Dragomir Anguelov, John Flynn, and Jana Kosecka. 3d bounding box estimation using deep learning and geometry. In *CVPR*, pages 7074–7082, 2017. 3

[32] Dennis Park, Rares Ambrus, Vitor Campanholo Guizilini, Jie Li, and Adrien Gaidon. Is pseudo-lidar needed for monocular 3d object detection? In *ICCV*, 2021. 7

[33] Liang Peng, Fei Liu, Zhengxu Yu, Senbo Yan, Dan Deng, and Deng Cai. Lidar point cloud guided monocular 3d object detection. *ArXiv*, abs/2104.09035, 2021. 2, 7

[34] Siyuan Qiao, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Vip-deeplab: Learning visual perception with depth-aware video panoptic segmentation. In *CVPR*, 2021. 3

[35] Zengyi Qin, Jinglu Wang, and Yan Lu. Monogrnet: A general framework for monocular 3d object detection. *TPAMI*, 2021. 7

[36] Cody Reading, Ali Harakeh, Julia Chae, and Steven L. Waslander. Categorical depth distribution network for monocular 3d object detection. In *CVPR*, pages 8551–8560, 2021. 3

[37] Cody Reading, Ali Harakeh, Julia Chae, and Steven L Waslander. Categorical depth distribution network for monocular 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8555–8564, 2021. 7

[38] Xuepeng Shi, Qianru Ye, Xiaozhi Chen, Chuangrong Chen, Zhixiang Chen, and Tae-Kyun Kim. Geometry-based distance decomposition for monocular 3d object detection. *ArXiv*, abs/2104.03775, 2021. 7

[39] Andrea Simonelli, Samuel Rota Bulò, Lorenzo Porzi, Peter Kontschieder, and Elisa Ricci. Are we missing confidence in pseudo-lidar methods for monocular 3d object detection? In *ICCV*, pages 3205–3213, 2021. 6

[40] Andrea Simonelli, Samuel Rota Bulò, Lorenzo Porzi, Manuel López-Antequera, and Peter Kontschieder. Disentangling monocular 3d object detection. *ICCV*, pages 1991–1999, 2019. 6

[41] Kihyuk Sohn, David Berthelot, Chun-Liang Li, Zizhao Zhang, Nicholas Carlini, Ekin Dogus Cubuk, Alexey Kurakin, Han Zhang, and Colin Raffel. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. In *NeurIPS*, 2020. 1, 3

[42] Kihyuk Sohn, Zizhao Zhang, Chun-Liang Li, Han Zhang, Chen-Yu Lee, and Tomas Pfister. A simple semi-supervised learning framework for object detection. In *arXiv:2005.04757*, 2020. 1, 3

[43] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, Jiquan Ngiam, Hang Zhao, Aleksei Timofeev, Scott M. Ettinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Yu Zhang, Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov. Scalability in perception for autonomous driving: Waymo open dataset. In *CVPR*, pages 2443–2451, 2020. 6

[44] Yihe Tang, Weifeng Chen, Yijun Luo, and Yuting Zhang. Humble teachers teach better students for semi-supervised object detection. In *CVPR*, 2021. 1, 3

[45] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *NeurIPS*, 2017. 2, 3

[46] Haiquan Wang, Yezhen Cong, Or Litany, Yue Gao, and Leonidas J. Guibas. 3dioumatch: Leveraging iou prediction for semi-supervised 3d object detection. In *CVPR*, 2021. 1, 3, 6

[47] Li Wang, Liang Du, Xiaoqing Ye, Yanwei Fu, Guodong Guo, X. Xue, Jianfeng Feng, and Li Zhang. Depth-conditioned dynamic message propagation for monocular 3d object detection. In *CVPR*, pages 454–463, 2021. 7

[48] Li Wang, Li Zhang, Yi Zhu, Zhi Zhang, Tong He, Mu Li, and X. Xue. Progressive coordinate transforms for monocular 3d object detection. In *NeurIPS*, 2021. 7

[49] Yan Wang, Wei-Lun Chao, Divyansh Garg, Bharath Hariharan, Mark Campbell, and Kilian Q Weinberger. Pseudo-lidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving. In *CVPR*, pages 8445–8453, 2019. 3

[50] Zhenyu Wang, Yali Li, Ye Guo, Lu Fang, and Shengjin Wang. Data-uncertainty guided multi-phase learning for semi-supervised object detection. In *CVPR*, 2021. 1, 3

[51] Qizhe Xie, Eduard H. Hovy, Minh-Thang Luong, and Quoc V. Le. Self-training with noisy student improves imagenet classification. In *CVPR*, pages 10684–10695, 2020. 1, 2, 3

[52] Mengde Xu, Zheng Zhang, Han Hu, Jianfeng Wang, Lijuan Wang, Fangyun Wei, Xiang Bai, and Zicheng Liu. End-to-end semi-supervised object detection with soft teacher. *ArXiv*, abs/2106.09018, 2021. 1, 3

[53] Qize Yang, Xihan Wei, Biao Wang, Xia Hua, and Lei Zhang. Interactive self-training with mean teachers for semi-supervised object detection. In *CVPR*, 2021. 1, 3

[54] Yurong You, Yan Wang, Wei-Lun Chao, Divyansh Garg, Geoff Pleiss, Bharath Hariharan, Mark Campbell, and Kilian Q Weinberger. Pseudo-lidar++: Accurate depth for 3d object detection in autonomous driving. In *ICLR*, 2020. 3

[55] Bowen Zhang, Yidong Wang, Wenxin Hou, Hao Wu, Jindong Wang, Manabu Okumura, and Takahiro Shinozaki. Flexmatch: Boosting semi-supervised learning with curriculum pseudo labeling. In *NeurIPS*, 2021. 1, 3

[56] Fangyuan Zhang, Tianxiang Pan, and Bin Wang. Semi-supervised object detection with adaptive class-rebalancing self-training. *ArXiv*, abs/2107.05031, 2021. 1, 3

[57] Renrui Zhang, Hang Qiu, Tai Wang, Xuan Xu, Ziyu Guo, Yu Jiao Qiao, Peng Gao, and Hongsheng Li. Monodetr: Depth-aware transformer for monocular 3d object detection. In *CVPR*, 2022. 7

[58] Yunpeng Zhang, Jiwen Lu, and Jie Zhou. Objects are different: Flexible monocular 3d object detection. In *CVPR*, pages 3288–3297, 2021. 2, 3, 6, 7, 8, 9

[59] Na Zhao, Tat-Seng Chua, and Gim Hee Lee. Sess: Self-ensembling semi-supervised 3d object detection. In *CVPR*, pages 11076–11084, 2020. 1, 3

[60] Wu Zheng, Weiliang Tang, Li Jiang, and Chi-Wing Fu. Se-ssd: Self-ensembling single-stage object detector from point cloud. In *CVPR*, pages 14489–14498, 2021. 3

[61] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *ArXiv*, abs/1904.07850, 2019. 3

[62] Yunsong Zhou, Yuan He, Hongzi Zhu, Cheng Wang, Hongyang Li, and Qinhong Jiang. Monocular 3d object detection: An extrinsic parameter free approach. In *CVPR*, pages 7552–7562, 2021. 7

[63] Zhikang Zou, Xiaoqing Ye, Liang Du, Xianhui Cheng, Xiao Tan, Li Zhang, Jianfeng Feng, X. Xue, and Errui Ding. The devil is in the task: Exploiting reciprocal appearance-localization features for monocular 3d object detection. In *ICCV*, pages 2693–2702, 2021. 7