# 3D Data Augmentation for Driving Scenes on Camera

Wenwen Tong[1*]   Jiangwei Xie[1*]   Tianyu Li[2*]   Hanming Deng[1*]   Xiangwei Geng[2]
Ruoyi Zhou[1]   Dingchen Yang[2]   Bo Dai[2]   Lewei Lu[1]   Hongyang Li[2✉]

[1]SenseTime Research    [2]Shanghai AI Laboratory

## Abstract

*Driving scenes are extremely diverse and complicated that it is impossible to collect all cases with human effort alone. While data augmentation is an effective technique to enrich the training data, existing methods for camera data in autonomous driving applications are confined to the 2D image plane, which may not optimally increase data diversity in 3D real-world scenarios. To this end, we propose a 3D data augmentation approach termed Drive-3DAug, aiming at augmenting the driving scenes on camera in the 3D space. We first utilize Neural Radiance Field (NeRF) to reconstruct the 3D models of background and foreground objects. Then, augmented driving scenes can be obtained by placing the 3D objects with adapted location and orientation at the pre-defined valid region of backgrounds. As such, the training database could be effectively scaled up. However, the 3D object modeling is constrained to the image quality and the limited viewpoints. To overcome these problems, we modify the original NeRF by introducing a geometric rectified loss and a symmetric-aware training strategy. We evaluate our method for the camera-only monocular 3D detection task on the Waymo and nuScences datasets. The proposed data augmentation approach contributes to a gain of 1.7% and 1.4% in terms of detection accuracy, on Waymo and nuScences respectively. Furthermore, the constructed 3D models serve as digital driving assets and could be recycled for different detectors or other 3D perception tasks.*
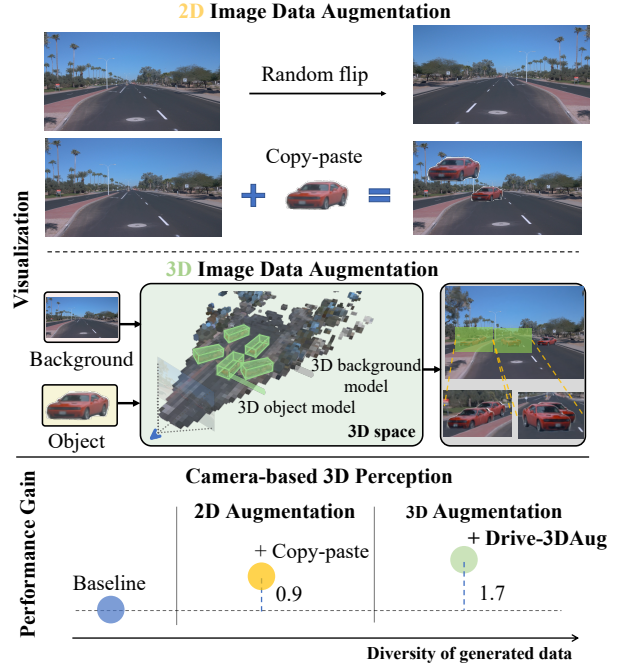
Figure 1. **Visualization of different data augmentation techniques and their *performance gain* of LET-AP [19] for monocular 3D detection on Waymo [42]** with 100 scenes augmented. Compared with previous 2D data augmentations, our Drive-3DAug method, modifies the driving scenes in the 3D space. This can generate more diverse driving scenes and contribute larger performance gain to camera-based 3D perception tasks. Baseline is the FCOS3D [45] method.

## 1. Introduction

3D perception system, particularly 3D object detection, plays a vital role for autonomous driving. Despite recent progress [1, 22, 50, 23, 28, 34, 35, 39, 54, 55, 4, 32, 45, 47, 48], the current perception system still suffers from the hard case challenge due to the long-tail driving scenes on the road, *e.g.* trucks with diversified pose on the road. To over-

come this challenge, data augmentation has been proven to be an effective technique to enrich training data. For LiDAR-based 3D perception, different data augmentation methods [36, 11, 16] have made great achievements by generating new drive scenes. However, it is still under-explored how to augment the driving scenes for camera-based 3D perception with data augmentation.

As illustrated in Figure 1, existing image data augmentation approaches are mostly restricted to the 2D image plane, such as image transformations [7] and copy-paste [14, 25].

---

*: Equal contribution.
✉: Corresponding author at lihongyang@pjlab.org.cn.

These techniques face challenges in changing the view of the components in the scene, such as rotating objects within the image, thereby limiting the diversity of generated driving scenes. By contrast, the data augmentation approaches [11, 16, 49] for point clouds are applied in the 3D space, offering more degrees of freedom to change the driving scenes. Although existing data augmentation approaches for image data have achieved some performance gains for the 3D detection task, such gains are still limited compared with the improvement brought by the 3D data augmentation approaches [36] for the LiDAR-based 3D perception tasks. This means the diversity of generated scenes is essential to improve the performance of 3D perception tasks. The manner of augmenting data in the 3D space is an effective way to create diverse scenes. One might argue that simulator [10, 38] is a powerful tool to generate synthetic 3D imagery to supplement the database. However, the sim2real bottleneck is a long-standing issue to be resolved.

In this work, we have pioneered research into 3D data augmentation for camera-based 3D percpetion for in autonomous driving. To implement 3D data augmentation for image data, we need to convert the scenes to the 3D space. This is because manipulating objects on 2D image plane satisfying 3D-imaging constraint, such as rotating pixels of objects, will generate flawed images. One desirable solution is using Neural Radiance Field (NeRF) [27, 33, 12] to reconstruct the 3D models of background and foreground objects obtained by decomposing the scene, then we can compose them for data augmentation. However, it is hard to achieve perfect decomposition without ground truth to extract objects. Pixels of the object will be mixed with the background pixels near the object edge. In addition, the limited viewpoints of objects in driving scenes make it difficult to apply NeRF for generating objects of novel views by large rotation, further limiting the diversity of generated scenes.

To this end, we present a novel 3D data augmentation approach for 2D images, named Drive-3DAug. Our approach has two stages. The first stage is to build the 3D models. We decompose the driving scenes into multiple backgrounds and foreground objects and use voxel-based NeRF [41] to turn them into 3D models. To overcome the difficulties of reconstructing the driving scenes, we propose a geometric rectified loss to weaken the effect of noisy edges and a symmetric-aware training strategy for objects with symmetry, such as vehicles, to broaden the available rendering viewpoints. After the first stage, we re-compose the 3D models of backgrounds and foreground objects to create new driving scenes. Considering the physical constraints in the real world, we adopt a strategy to identify the valid region of the backgrounds. The augmented images are then generated by rendering the constructed scenes and can be further applied to the training of 3D perception tasks. Com-

pared with existing approaches[25], our method contributes to a larger performance gain of the 3D detection task, as described in Figure 1.

We evaluate Drive-3DAug for the monocular 3D detection task, as it is one of the most important 3D perception tasks in autonomous driving, on the Waymo [42] and nuScenes [3] datasets with different detectors. Our method is able to achieve 1.7% improvement of performance on Waymo, especially 2.2% for vehicles with rare orientations, and 1.4% on nuScenes on their corresponding metrics. Moreover, once the 3D models of the backgrounds and foreground objects are reconstructed, they can serve as a digital driving asset, which can be used repeatedly for different tasks.

To sum up, our contributions are three-fold:

- We have pioneered research into the 3D data augmentation problem for camera-based 3D perception in autonomous driving.

- We propose a 3D data augmentation approach based on NeRF including improvements by a geometrically rectified loss and a symmetric-aware training strategy to generate more natural and diverse images.

- We evaluate our method on the Waymo and nuScenes dataset, demonstrating that it can improve the performance of 3D object detection, and the 3D models can be further used as digital driving assets.

## 2. Related Work

**Data Augmentation for 3D perception.** Data augmentation is a powerful technique to improve performances of perception algorithms [40, 56]. For 3D perception, augmentation methods vary in data modalities. For point cloud data, flipping, rotation and translation of objects and backgrounds are common techniques [5, 17, 51, 6]. For image data, 2D augmentation methods can be lifted to 3D with geometry constraints. [25] improve random scale, crop and copy-paste from 2D to 3D with 2D–3D geometry relationship. For multi-modality data, [53, 44] are proposed to keep a consistency between images and point clouds. Although these methods improve the performance of 3D perception tasks, the operations on the image plane are much less flexible than the 3D data augmentation for point clouds, *e.g.*, unable to rotate vehicles. Besides, aggressive augmentation techniques on images always violate geometry constraints and lead to unnatural, flawed data. In comparison, the approach we proposed can generate images of diverse driving scenes with more degrees of freedom by object manipulation in the 3D space, which takes geometry constraints and occlusions into account at the same time.
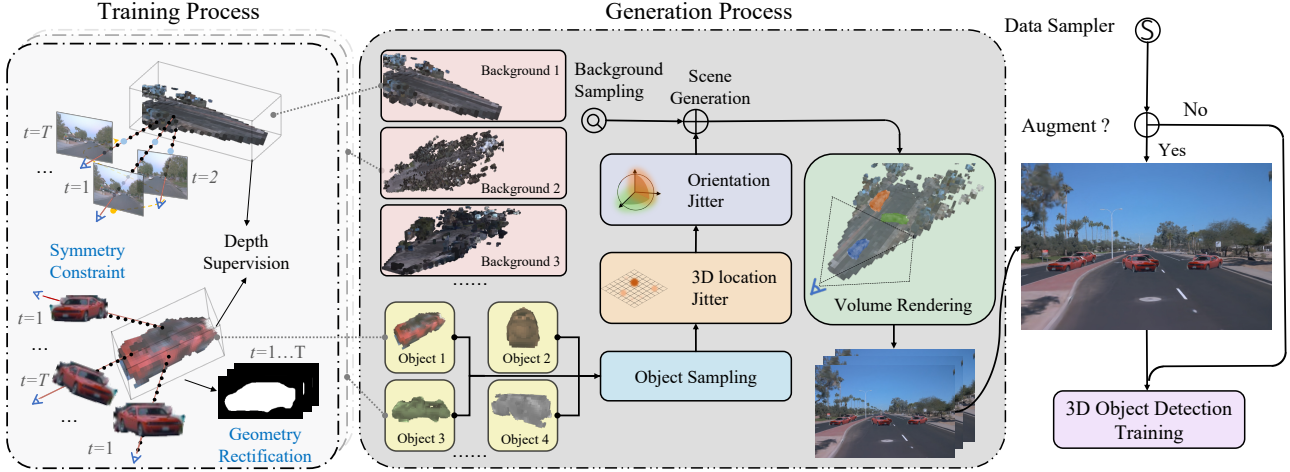
Figure 2. **Overview of Drive-3DAug** for 3D data augmentation. Driving scenes are decomposed into multiple backgrounds and objects. For each background and object, we use multi-frame views to reconstruct them separately by voxel-based NeRF [41]. To further improve the reconstruction quality, we introduce the symmetry constraint, geometry rectification, and depth supervision for the NeRF. We edit the scene in 3D space with the manipulation of trained 3D models, and the images are generated by rendering the composed new scenes for the following 3D perception tasks.

**NeRF for Scene Generation.** NeRF [27] is a powerful tool for novel view synthesis, which represents a scene with a fully connected neural network and optimizes it with differentiable volume rendering. It has been recently applied to autonomous driving scenarios [43, 24, 29, 13]. Block–NeRF [43] reconstructs a whole city by merging multiple block-NeRFs with predicted visibility. Auto-RF [29] focuses on vehicle reconstruction in autonomous driving scenes. Considering the majority of NeRF methods are per-scene fitting, how to quickly reconstruct a scene is an important challenge because of the large scale of driving scenes. To overcome this, voxel-based NeRF [12, 41, 30] and depth-supervised NeRF [9], which are adopted in our method, have been proposed to accelerate the training of NeRF. Moreover, some works begin to explore how to edit scenes with NeRF [31, 21, 24] for autonomous driving. However, these works are not extended to augmenting the data for 3D perception tasks. PNF [21] does not consider driving data only has few views for most objects, which restricts the quality of novel view synthesis. READ [24] realizes autonomous driving scene editing with larger viewpoints of objects by point-based NeRF. However, this method relies on LiDAR data which limits its application scope. By contrast, the LiDAR data is not necessary for our method.

## 3. Method

Our goal is to create diverse driving scenes for improving 3D perception tasks, especially for camera-based 3D detection, by aid of 3D data augmentation. To this end, we propose Drive-3DAug.

### 3.1. Overview

As demonstrated in Figure 2, Drive-3DAug has two stages for implementing 3D data augmentation.

**Stage 1 - Training.** The first stage is to construct the 3D models from images for the following generation of new driving scenes. To achieve this, we decompose the scenes into backgrounds and foregrounds and build the 3D models of them by training the NeRF. Then we can edit the scenes in the 3D space. Considering the characteristics of driving data, we further develop several techniques to improve the training process.

**Stage 2 - Generation.** The second stage is to augment the training data through the 3D models. To create new driving scenes, we combine the models of foreground objects and background scenes in the 3D space with the manipulation of objects including the location and orientation jittering. In order to make the generated scenes close to the real-world driving scene, e.g., we cannot place vehicles on a tree, we design a strategy to identify the valid region of the background scene where we place the foreground object. Then we use volume rendering to generate new images of these scenes for the training of camera-based 3D perception tasks.

### 3.2. 3D Model Training

To edit driving scenes for 3D data augmentation, we need to construct a set of 3D models for backgrounds and foregrounds. To achieve this, we first utilize an off-the-shelf instance segmentation model to extract the objects from the

backgrounds. Then we use the Intersection of Union (IoU) of the object masks and the projection of 3D boxes on images to match the extracted object and the 3D annotation. After this, we use NeRF to reconstruct the 3D models from images. We reconstruct the 3D objects based on the extracted masks in consecutive frames by matching the object masks with IoU constraint, and we only consider the totally visible objects with intact masks for 3D reconstruction. To model the static backgrounds, the moving object masks in the images are filtered.

To efficiently reconstruct the backgrounds and objects, we use voxel-based NeRF [12, 41, 30] instead of MLP-based NeRF [27, 33]. Voxel-based NeRF has a density voxel grid $V_{density}$, and a feature voxel grid $V_{color}$ with shallow MLPs to represent the geometry and appearance respectively. The NeRF model is trained by minimizing the loss between the rendered pixel color $C(r)$ and observed pixel color $\hat{C}(r)$ along the ray $r$ given by

$$\mathcal{L}_{Color} = \sum_{r \in \mathcal{R}(\mathbf{P})} \|\hat{C}(r) - C(r)\|_2^2, \qquad (1)$$

where $\mathcal{R}(\mathbf{P})$ is the set of rendered rays in a batch. To accelerate the model training, we introduce depth supervision to optimize the voxel field. Then the NeRF model is trained by minimizing the loss

$$\mathcal{L} = \mathcal{L}_{color} + \mathcal{L}_{Depth}, \qquad (2)$$

where $\mathcal{L}_{Depth}$ is the $L_1$ loss between the rendered depth and observed depth. We can use LiDAR data or estimated depth by methods such as as SFM [37] or PACKNet [15] to obtain the observed depth. We adopt the LiDAR data in this work as it is common in driving data. For the convenience of scene editing, we reconstruct the static background of the scene in the world coordinate system and the object modeling the local object 3D box coordinate system. Once we finish the training of the background and object NeRFs, they can serve as the digital driving assets for the repeated generation of novel scenes.

**Analysis.** Although the 3D models trained through this progress can be used for novel scene creation, they still struggle with several problems. First, because of the complexity of the driving scenes and the imperfect instance segmentation model, the extracted images of objects usually suffer from edge defects. NeRF can not model the object geometry well based on the noise mask as the background pixel may leak into the object model. In addition, the voxel grid representation through interpolation in voxel-based NeRF can further hinder the modeling of clear geometry near object boundaries. Second, the range of available viewpoints for objects is very important for creating diverse driving scenes. However, the objects in the driving data
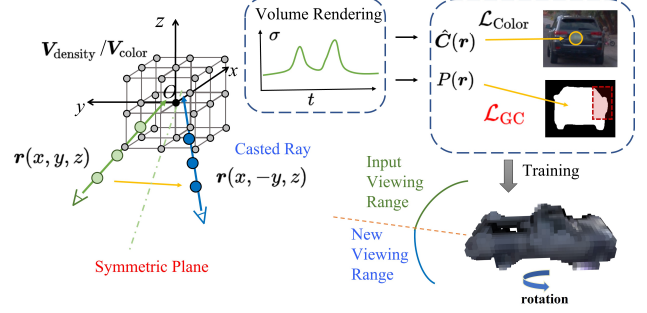


Figure 3. **Modified NeRF** including a geometric rectified loss and a symmetric-aware training strategy. They can alleviate the effects of imperfect object extraction and increase the range of viewpoints, respectively.

have limited views, indicating that we can only render objects within a small degree. This will restrict the diversity of the generated driving scenes. In terms of this, we refine the voxel-based NeRF model to promote the object models.

### 3.3. Improvements of Object Models

We design a geometric rectified loss and a symmetric-aware training strategy for the object models to ensure the quality and diversity of the generated novel scenes.

**Geometric Rectified Loss.** To avoid the effects of imperfect object extraction, we propose a geometric rectified loss to correct the geometry of the object model. As illustrated in Figure 3, we add an auxiliary task for the training of the NeRF model by the classification of the pixel as the foreground or background pixel. Because the edge defects are different across consecutive frames, the temporal inconsistency can make the model remove the edge defects. Specifically, the probability of a rendered pixel being an object pixel can be approximated by

$$P(r) = 1 - \exp\left(-\int_{t_a}^{t_b} \sigma ds\right), \qquad (3)$$

where $t_a$ and $t_b$ denote the entrance and exit points of the ray-object intersection respectively, and $\sigma$ is the density of sampled point in the ray. Then we implement the geometric rectified loss as

$$\mathcal{L}_{GC} = - \sum_{r \in \mathcal{R}(\mathbf{P})} log P(r), \qquad (4)$$

to decrease the voxel density near the mask edge, where $\mathcal{R}(\mathbf{P})$ is the set of rendered rays in a batch. The final loss for training the object voxel field is defined by

$$\mathcal{L}_{object} = \mathcal{L}_{Color} + \mathcal{L}_{Depth} + \mathcal{L}_{GC}. \qquad (5)$$
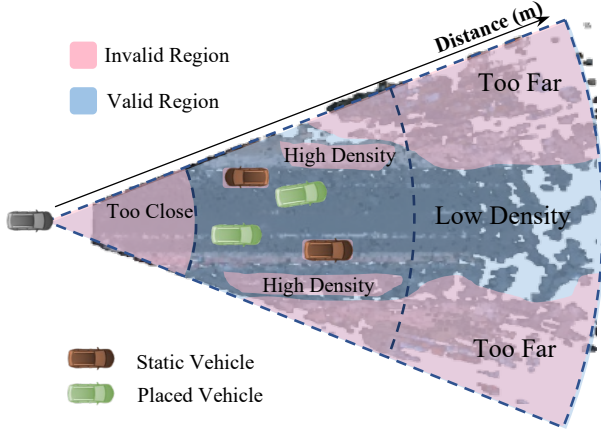
Figure 4. **Valid region** identification for the placement of vehicles, when composing a new driving scene.

**Symmetric-aware Training Strategy.** To enrich the viewpoints of objects, we design a symmetric-aware training strategy. Considering that objects, such as vehicles, are usually geometric symmetric in the driving scenes, we can create a symmetric object voxel field for them. As depicted in Figure 3, given a pixel $p_r$ along the ray $r(x, y, z)$, we can create a symmetric virtual camera that casts a ray $r(x, -y, z)$ and targets the pixel $p_l$, in which $p_r = p_l$, supposing the symmetry of the object. In this way, we increase the number of camera views for the object model training. Then we can render novel views of the object rotated with a larger degree.

### 3.4. Generation of Augmented Driving Scenes

As indicated in Figure 2, we generate a new driving scene through the combination of the object models and background models. We can place objects following the original location and orientation distributions of objects in data or change the distributions, such as making samples more balanced. Besides, to make the placement of objects satisfy the physical law or other constraints in the real world, such as car on the road instead of the sky, we also identify the valid region of the background before the placement. Otherwise, the models of 3D perception tasks cannot learn the proper context information. As shown in Figure 4, we define the valid region on the bird's eye view based on the density field of the background model, which is inspired by the Lidar-Aug [11]. To find the valid region, we first divide the 3D space into sets of pillars, and then the region can be split into valid and invalid states based on the density distribution of voxel in the corresponding pillar. Specifically, the valid region satisfies the low-density constraints:

$$\max(Z_p) < \delta_1, \text{ and } mean(Z_p) < \delta_2, \quad (6)$$

where $Z_p$ is the array denoting the density of points in the pillar, and $\delta_1$ and $\delta_2$ are hyper-parameters. This represents

no large object is in this region. Furthermore, we can filter the low-density region behind the high-density region such as the wall. In this way, we can put the object in an appropriate position. To avoid the collision, we calculate 3D IoUs between the placed objects and existing objects to ensure consistency between the foreground objects and the backgrounds. Then we jointly render the merged 3D model to generate images of the new scenes. The generated images can be directly applied for the downstream 3D detection task to improve the performance of detectors.

## 4. Experiments

We validate the proposed Drive-3DAug method on the monocular 3D detection task, which is one of the most important 3D perception tasks in autonomous driving.

### 4.1. Datasets and Metrics

**Waymo Dataset.** The Waymo Open Dataset [42] is a large-scale dataset for autonomous driving that contains 798 scenes in the training dataset and 202 scenes in the validation dataset. The image resolution for the front camera is $1920 \times 1280$. Waymo uses the LET-AP [20], the average precision with longitudinal error tolerance, to evaluate detection models. Besides, Waymo also adopts the LET-APL and LET-APH metrics, which are the longitudinal affinity weighted LET-AP and the heading accuracy weighted LET-AP, respectively.

**nuScenes Dataset.** The nuScenes [3] is a widely used benchmark for 3D object detection. It contains 700 training scenes and 150 validation scenes. The resolution of each image is $1600 \times 900$. As for the metrics, nuScenes computes mAP using the center distance on the ground plane to match the predicted boxe and the ground truth. It also contains different types of true positive metrics (TP metrics). We use ATE, ASE and AOE in this paper, for measuring the errors of translation, scale and orientation, respectively.

### 4.2. Implementation Details

**Building Digital Driving Asset.** We use the SOLO v2 [46] trained on COCO as the instance segmentation model for scene decomposition. We consider *vehicle* (Waymo) or *car* (nuScenes) as the foreground objects, since they are the most important components in driving scenes For 3D model reconstruction, we use the same model configuration as the DVGO [41] with our proposed techniques. We use 30-40 consecutive frames spanning an area of about 100-200 meters as one background and train each background model with 40,000 iterations. For the background voxel grid, we set the resolution as $330^3$ with a voxel size of 0.25-0.3m. The object models are trained with 20-60 consecutive frames, and we set the voxel size as 0.25m consistent with the background voxel size. The grid point number

| Method | LET-AP | LET-APH | LET-APL | LET-APL [50m, +∞) | LET-APH ∼45° | LET-APH ∼90° | LET-APH ∼135° |
|---|---|---|---|---|---|---|---|
| FCOS3D [45] | 0.585 | 0.573 | 0.393 | 0.278 | 0.293 | 0.285 | 0.223 |
| + Copy-paste [25] | 0.594 | 0.581 | 0.401 | 0.290 | 0.295 | 0.283 | **0.234** |
| + Drive-3DAug w/o RT | 0.595 | 0.583 | 0.404 | 0.287 | 0.293 | 0.295 | 0.226 |
| + Drive-3DAug w/ RT | **0.602** | **0.590** | **0.410** | **0.298** | **0.315** | **0.299** | **0.234** |
| SMOKE [26] | 0.586 | 0.579 | 0.417 | 0.304 | 0.312 | 0.291 | 0.264 |
| + Copy-paste [25] | 0.594 | 0.587 | 0.425 | 0.320 | 0.303 | 0.301 | 0.239 |
| + Drive-3DAug w/o RT | 0.592 | 0.584 | 0.421 | 0.318 | 0.314 | 0.298 | **0.266** |
| + Drive-3DAug w/ RT | **0.598** | **0.590** | **0.426** | **0.322** | **0.333** | **0.303** | 0.255 |

Table 1. **Monocular 3D detection results** on the Waymo validation set. The left part shows the main metrics. LET-AP represents longitudinal error tolerant 3D average precision [20]. LET-APH and LET-APL represent LET-AP penalized by heading errors and longitudinal localization errors, respectively. The right part of the table shows metrics under specific hard cases. The [50m, +∞) metric only calculates AP over objects with distance to ego-car greater than 50m. The ∼h° metric only consider the objects with |heading| near. It is observed that our *Drive-3DAug w/ RT* performs consistently better than any other setting.



(a) Reference

(b) Placing the object randomly

(c) Placing the object at a distance

(d) Generating the dense scene

Figure 5. **Scene generation** with different placement strategies of the *vehicle*.



Figure 6. **Distribution of heading direction** for *vehicle* on the Waymo training set and the distribution augmented by Drive-3DAug.

is about 1,000. Considering the construction cost and limitation of NeRF model for extreme illumination conditions, we select a subset of 100 sunny scenarios for each dataset. We use them for the data augmentation.

**Applying Data Augmentation.** Our method generates new data through rendering the recomposed scenes of the randomly selected 3D models. We manipulate the object in 3D space by 3D location jittering and orientation jittering and place it on the valid region of arbitrary backgrounds. The location jittering is defined by the maximum translation $(T_x, T_y)$ along the x-direction and the y-direction, and the orientation jittering is given by maximum rotation angle $T_\theta$. We consider two data augmentation strategies to valid the effectiveness of adding rotation and translation for 3D perception. The first is *Drive-3DAug w/o RT*, in which we set translation $T_x = 0$, $T_y = 0$, and rotation $T_\theta = 0$ and
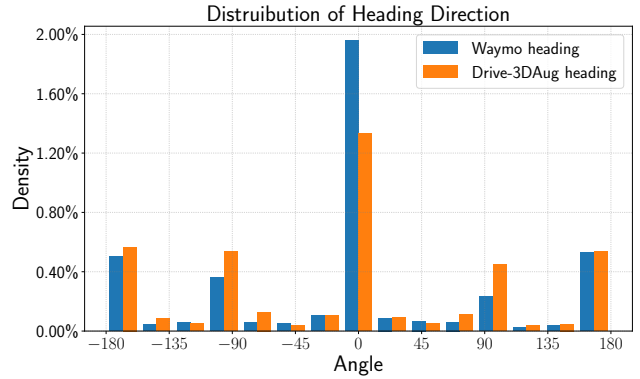
1-2 new objects are arbitrarily pasted into the background scene on average. The second is *Drive-3DAug w/ RT* with $T_x = 20m$, $T_y = 5m$, and $T_\theta = 30°$. As for determining the valid region in the background scene, we set the pillar $Z_p$ resolution as 2m×2m, $\delta_1 = 30$, and $\delta_2 = 15$ in Eq. 6. We generate 12 new images for every background model. This progress is offline data augmentation and these images are repeatedly used by different detectors.

**Training Detectors.** Two typical camera-based monocular 3D object detectors, FCOS3D[45] and SMOKE[26] are utilized to investigate the performance of our proposed method as they are two of the most popular and commonly used mono3D detectors. We maintain the same hyperparameters of detectors for all experiments in our study, introduced in the supplementary. We use all the scenes in the training set to train the detectors and evaluate the detectors on the entire validation set. However, we sample

| Method | AP | ATE | ASE | AOE |
|---|---|---|---|---|
| FCOS3D[45] | 0.319 | 0.739 | 0.160 | 0.096 |
| + Copy-paste [25] | 0.324 | 0.721 | **0.156** | 0.123 |
| + Drive-3DAug w/ RT | **0.333** | **0.705** | 0.158 | **0.092** |

Table 2. **Monocular 3D detection results** on the nuScenes validation set. ATE, ASE, and AOE are used for measuring translation, scale, and orientation errors respectively.



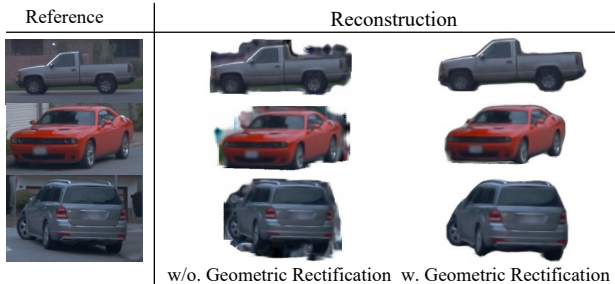w/o. Geometric Rectification    w. Geometric Rectification

Figure 7. **Visualization** of the ablative study for the geometric rectified loss. After applying such a loss, the reconstructed object models have much fewer edge defects.

| GRL | SAT | LET-AP | LET-APH | LET-APL |
|---|---|---|---|---|
| - | - | 0.590 | 0.578 | 0.403 |
| ✓ | - | 0.596 | 0.583 | 0.407 |
| ✓ | ✓ | **0.602** | **0.590** | **0.410** |

Table 3. **Ablation Study** of Drive-3DAug w R/T for FCOS3D on Waymo validation set. GRL means the geometric rectified loss and SAT means symmetric-aware training.

training data every 3 frames on Waymo during training due to limited computational resources. Besides, although our method can be applied to any view of cameras, we only take into account images taken by the front camera on Waymo and nuScenes in this work because of the computational resource. For the usage of the generated images, we randomly replace the image in a batch with the augmented data if it belongs to the scenes in the digital driving asset.

## 4.3. Main Results

**Novel Scene Generation.** Drive-3DAug can generate diverse scenes with the ability to manipulate objects in the 3D space. As indicated in Figure 5, our method can place the car at any distance or generate a very dense scene to augment the data. Moreover, previous image data augmentation methods are likely to put the car on the sky because they do not define the valid region. By contrast, our placement obeys the real world situation. Besides the placement, our method can also alleviate the imbalance problem of vehicle orientation with orientation jittering to objects. Figure 6 shows the extremely imbalanced distribution of vehicle heading direction on the Waymo dataset, where most heading directions are at $0°$ or $180°$. In comparison, the distribution of vehicle heading distribution is more balanced than the original distribution, although the augmented distribution is still imbalanced. This is because of the randomly jitter for the orientation of objects.

**3D Object Dectection.** Table 1 illustrates the main results on the Waymo validation set. We compare our Drive-

3DAug with the geometry-consistent copy-paste method [53], which copies and pastes objects across different scenes with depth constraints and camera transformation of 3D annotations. Besides, we also apply the valid region strategy for this method with the help of LiDAR data. It also uses the selected 100 scenes for augmentation. Since our data augmentation method is applied only to vehicles, we use the LET-AP on *vehicle* to evaluate the effectiveness of our approach. Compared to *baseline*, our *NeRF w/ RT* approach significantly improves LET-AP by 1.5%, LET-APH by 1.4%, and LET-APL by 1.5% on average across both FCOS3D and SMOKE. The consistent improvement between these detectors indicates that our approach is robust to different backbone and head architectures. It also outperforms *Copy-paste* by 0.6%, 0.6%, and 0.5% on LET-AP, LET-APH, and LET-APL, respectively. Furthermore, we also compare the performance of different methods under far distances and different orientation ranges of vehicles. The right of Table 1 shows that the detector trained with *Drive-3DAug w/ RT* outperforms other methods under all the scenarios, especially for the objects with $\sim 45°$ orientation, which is the fewest in the original data. This means the ability to add larger translation and rotation to placed objects is essential for the data augmentation techniques in 3D perception.

Table 2 reports the experimental results on the nuScenes validation set. We can see *Drive-3DAug w/ RT* also improves the performance of class *car* on the nuScenes dataset. Additionally, it is worth noting that Drive-3DAug obtains a 1.6% improvement on ATE, a 3.1% improvement on AOE for the FCOS3D detector compared with the geometry-consistent copy-paste method. The result on ASE is comparable for these two methods. This is because we do not apply the scale transformation. This further proves that our method is capable of resolving the issues of previous methods with only 2D pixel movements that lack the ability to generate realistic rotated or translated objects.

## 4.4. Ablative Studies

Table 3 reports the results of our ablation experiments with the FCOS3D detector. The baseline is using DVGO [41] with depth supervision for data augmentation. The results indicate that even with vanilla voxel-based NeRF with-

| Rotation | **−10°** | **0°** | **10°** |
|---|---|---|---|
| w/o. symmetric training | | | |
| w. symmetric training | | | |

Figure 8. **Visualization** of the ablative study for the symmetric-aware training strategy. It can make the model generate more realistic images with larger degrees of rotation.

(a)

(b)

(c)

(d)

Figure 9. **Corner Cases** generated by the Drive-3DAug, where (a), (b) illustrate occlusion situations, (c) demonstrates strange vehicle heading direction, and (d) shows a vehicle appended on slope road. Cars with 3D bounding boxes are rendered by NeRF.

out any improvement, our proposed Drive-3DAug can improve the performance of the detector by 0.5% on LET-AP. After adding the geometric rectified loss and the symmetric-aware training strategy, the LET-AP is further increased by 0.6% and 0.6% respectively. In addition, we provide some visualizations to demonstrate the effectiveness of each component. Figure 7 shows that the geometric rectified loss can significantly suppress the edge defects and improve the reconstruction quality. Figure 8 shows the results of the symmetric-aware training strategy for the novel view synthesis. We can see that adding this training strategy makes the model have larger viewpoints.

**Reconstruction Cost.** Table 4 depicts the comparison of reconstruction speed on a NVIDIA V100 GPU between previous methods and our method. Unlike the MLP-based NeRF [27] which needs more than 20 hours of training for one background, it takes about 0.5h for the voxel-based NeRF with depth supervision. For the object model, the reconstruction time is within minutes. Considering the model size of our NeRF is rather small, we can run multiple reconstructions in parallel. Moreover, once these models are trained, they could be recycled for different detectors as a

| Method | Voxel Field. | Depth Sup. | Cost |
|---|---|---|---|
| NeRF [27] | - | - | > 20h |
| DVGO [41] | ✓ | - | 0.7h |
| Ours | ✓ | ✓ | 0.5**h** |

Table 4. **Reconstruction cost** of different methods for one background, assessed on a NVIDIA V100 GPU with 16GB memory.

| Method | LET-AP | |
|---|---|---|
| | Pedestrian | Cyclist |
| FCOS3D [45] | 0.339 | 0.231 |
| + Drive-3DAug w RT | **0.344** | **0.234** |

Table 5. **Monocular 3D detection results** of pedestrian and cyclist on Waymo validation set. LET-AP is the longitudinal error tolerant 3D average precision.

digital driving assets.

**Corner Case Generation.** Drive-3DAug is able to generate many photographic data for various corner cases [2] without much effort for autonomous driving systems. As described in Figure 9, we use our method to simulate several corner cases including *the car occluded by the environment*, *the car appearing on the road with strange positions and headings*, and *the car on the slope*, which are hard to collect in the real world. This shows that 3D data augmentation can help alleviate issues of autonomous driving caused by plenty of corner cases.

**Pedestrian and Cyclist Augmentation.** Table 5 shows the results for applying Drive-3DAug to the *pedestrian* and *cyclist* on Waymo. These objects are difficult for reconstruction, especially in driving scenes. Because they are not rigid objects and may change their pose dramatically in consecutive frames. However, our method can still improve the detection results for these two classes by 0.5% and 0.3% on LET-AP, respectively.

## 5. Conclusion

In this paper, we propose Drive-3DAug, the first 3D data augmentation technique for camera-based 3D perception task in autonomous driving. We represent the scene with background and foreground 3D models by NeRF and randomly combine them to generate new driving scenes. The quality and diversity issues of generated scenes are addressed by our novel geometric rectified loss and symmetry-aware training strategies. We demonstrate the effectiveness of our method on multiple datasets and detectors. Furthermore, these 3D models can be regarded as the digital driving asset and benefit the community of this area.

**Discussion.** Driving scenes are extremely complicated, including lots of object categories under different illumination and weather conditions. Currently, our method only augments limited classes of objects under good illumination conditions. It is worth to include more situations as the digital driving asset for the future work.

# References

[1] Garrick Brazil and Xiaoming Liu. M3d-rpn: Monocular 3d region proposal network for object detection. In *ICCV*, pages 9287–9296, 2019. 1

[2] Jasmin Breitenstein, Jan-Aike Termöhlen, Daniel Lipinski, and Tim Fingscheidt. Corner cases for visual perception in automated driving: Some guidance on detection approaches. *arXiv preprint arXiv:2102.05897*, 2021. 8

[3] 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*, pages 11621–11631, 2020. 2, 5

[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. *NeurIPS*, 28, 2015. 1

[5] Shuyang Cheng, Zhaoqi Leng, Ekin Dogus Cubuk, Barret Zoph, Chunyan Bai, Jiquan Ngiam, Yang Song, Benjamin Caine, Vijay Vasudevan, Congcong Li, et al. Improving 3d object detection through progressive population based augmentation. In *ECCV*, pages 279–294, 2020. 2

[6] Jaeseok Choi, Yeji Song, and Nojun Kwak. Part-aware data augmentation for 3d object detection in point cloud. In *IROS*, pages 3391–3397, 2021. 2

[7] Ekin Dogus Cubuk, Barret Zoph, Dandelion Mané, Vijay Vasudevan, and Quoc V. Le. Autoaugment: Learning augmentation policies from data. *CoRR*, abs/1805.09501, 2018. 1

[8] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *ICCV*, pages 764–773, 2017. 11

[9] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *CVPR*, pages 12882–12891, 2022. 3

[10] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. CARLA: An open urban driving simulator. In *CORL*, pages 1–16, 2017. 2

[11] Jin Fang, Xinxin Zuo, Dingfu Zhou, Shengze Jin, Sen Wang, and Liangjun Zhang. Lidar-aug: A general rendering-based augmentation framework for 3d object detection. In *CVPR*, pages 4710–4720, 2021. 1, 2, 5

[12] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, pages 5501–5510, 2022. 2, 3, 4

[13] Xiao Fu, Shangzhan Zhang, Tianrun Chen, Yichong Lu, Lanyun Zhu, Xiaowei Zhou, Andreas Geiger, and Yiyi Liao. Panoptic nerf: 3d-to-2d label transfer for panoptic urban scene segmentation. In *3DV*, 2022. 3

[14] Golnaz Ghiasi, Yin Cui, Aravind Srinivas, Rui Qian, Tsung-Yi Lin, Ekin D Cubuk, Quoc V Le, and Barret Zoph. Simple copy-paste is a strong data augmentation method for instance segmentation. In *CVPR*, pages 2918–2928, 2021. 1

[15] Vitor Guizilini, Rares Ambrus, Sudeep Pillai, Allan Raventos, and Adrien Gaidon. 3d packing for self-supervised monocular depth estimation. In *CVPR*, pages 2485–2494, 2020. 4

[16] Martin Hahner, Dengxin Dai, Alexander Liniger, and Luc Van Gool. Quantifying data augmentation for lidar based 3d object detection. *CoRR*, abs/2004.01643, 2020. 1, 2

[17] Martin Hahner, Dengxin Dai, Alexander Liniger, and Luc Van Gool. Quantifying data augmentation for lidar based 3d object detection. *arXiv preprint arXiv:2004.01643*, 2020. 2

[18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 11

[19] Wei-Chih Hung, Henrik Kretzschmar, Vincent Casser, Jyh-Jing Hwang, and Dragomir Anguelov. Let-3d-ap: Longitudinal error tolerant 3d average precision for camera-only 3d detection. *arXiv preprint arXiv:2206.07705*, 2022. 1

[20] Wei-Chih Hung, Henrik Kretzschmar, Vincent Casser, Jyh-Jing Hwang, and Dragomir Anguelov. Let-3d-ap: Longitudinal error tolerant 3d average precision for camera-only 3d detection. *arXiv preprint arXiv:2206.07705*, 2022. 5, 6

[21] Abhijit Kundu, Kyle Genova, Xiaoqi Yin, Alireza Fathi, Caroline Pantofaru, Leonidas J. Guibas, Andrea Tagliasacchi, Frank Dellaert, and Thomas A. Funkhouser. Panoptic neural fields: A semantic object-aware neural scene representation. In *CVPR*, pages 12861–12871, 2022. 3

[22] 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. 1

[23] Peixuan Li, Huaici Zhao, Pengfei Liu, and Feidao Cao. Rtm3d: Real-time monocular 3d detection from object keypoints for autonomous driving. In *ECCV*, pages 644–660. Springer, 2020. 1

[24] Zhuopeng Li, Lu Li, Zeyu Ma, Ping Zhang, Junbo Chen, and Jianke Zhu. Read: Large-scale neural scene rendering for autonomous driving. *arXiv preprint arXiv:2205.05509*, 2022. 3

[25] Qing Lian, Botao Ye, Ruijia Xu, Weilong Yao, and Tong Zhang. Exploring geometric consistency for monocular 3d object detection. In *CVPR*, pages 1685–1694, 2022. 1, 2, 6, 7

[26] Zechen Liu, Zizhang Wu, and Roland Tóth. Smoke: Single-stage monocular 3d object detection via keypoint estimation. In *CVPRW*, pages 996–997, 2020. 6

[27] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, pages 99–106, 2020. 2, 3, 4, 8, 11

[28] Arsalan Mousavian, Dragomir Anguelov, John Flynn, and Jana Kosecka. 3d bounding box estimation using deep learning and geometry. In *CVPR*, pages 7074–7082, 2017. 1

[29] Norman Müller, Andrea Simonelli, Lorenzo Porzi, Samuel Rota Bulò, Matthias Nießner, and Peter Kontschieder. Autorf: Learning 3d object radiance fields from single view observations. In *CVPR*, pages 3971–3980, 2022. 3

[30] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, July 2022. 3, 4

[31] Julian Ost, Fahim Mannan, Nils Thuerey, Julian Knodt, and Felix Heide. Neural scene graphs for dynamic scenes. In *CVPR*, pages 2856–2865, 2021. 3

[32] Dennis Park, Rares Ambrus, Vitor Guizilini, Jie Li, and Adrien Gaidon. Is pseudo-lidar needed for monocular 3d object detection? In *ICCV*, pages 3142–3152, 2021. 1

[33] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *ICCV*, pages 5865–5874, 2021. 2, 4

[34] Zengyi Qin, Jinglu Wang, and Yan Lu. Monogrnet: A geometric reasoning network for monocular 3d object localization. In *AAAI*, volume 33, pages 8851–8858, 2019. 1

[35] Zengyi Qin, Jinglu Wang, and Yan Lu. Triangulation learning network: from monocular to stereo 3d object detection. In *CVPR*, pages 7615–7623, 2019. 1

[36] Matthias Reuse, Martin Simon, and Bernhard Sick. About the ambiguity of data augmentation for 3d object detection in autonomous driving. In *ICCVW*, pages 979–987, 2021. 1, 2

[37] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, pages 4104–4113, 2016. 4

[38] Shital Shah, Debadeepta Dey, Chris Lovett, and Ashish Kapoor. Airsim: High-fidelity visual and physical simulation for autonomous vehicles. In Marco Hutter and Roland Siegwart, editors, *Field and Service Robotics, Results of the 11th International Conference, FSR 2017, Zurich, Switzerland, 12-15 September 2017*, volume 5 of *Springer Proceedings in Advanced Robotics*, pages 621–635. Springer, 2017. 2

[39] Xuepeng Shi, Qi Ye, Xiaozhi Chen, Chuangrong Chen, Zhixiang Chen, and Tae-Kyun Kim. Geometry-based distance decomposition for monocular 3d object detection. In *ICCV*, pages 15172–15181, 2021. 1

[40] Connor Shorten and Taghi M Khoshgoftaar. A survey on image data augmentation for deep learning. *Journal of big data*, 6(1):1–48, 2019. 2

[41] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *CVPR*, pages 5459–5469, 2022. 2, 3, 4, 5, 7, 8, 11

[42] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *CVPR*, pages 2446–2454, 2020. 1, 2, 5

[43] Matthew Tancik, Vincent Casser, Xinchen Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. In *CVPR*, pages 8248–8258, 2022. 3

[44] Chunwei Wang, Chao Ma, Ming Zhu, and Xiaokang Yang. Pointaugmenting: Cross-modal augmentation for 3d object detection. In *CVPR*, pages 11794–11803, 2021. 2

[45] Tai Wang, Xinge Zhu, Jiangmiao Pang, and Dahua Lin. Fcos3d: Fully convolutional one-stage monocular 3d object detection. In *CVPR*, 2021. 1, 6, 7, 8

[46] Xinlong Wang, Tao Kong, Chunhua Shen, Yuning Jiang, and Lei Li. Solo: Segmenting objects by locations. In *ECCV*, pages 649–665, 2020. 5

[47] Yue Wang, Vitor Campagnolo Guizilini, Tianyuan Zhang, Yilun Wang, Hang Zhao, and Justin Solomon. Detr3d: 3d object detection from multi-view images via 3d-to-2d queries. In *CoRL*, pages 180–191. PMLR, 2022. 1

[48] Xinshuo Weng and Kris Kitani. Monocular 3d object detection with pseudo-lidar point cloud. In *ICCVW*, 2019. 1

[49] Aoran Xiao, Jiaxing Huang, Dayan Guan, Kaiwen Cui, Shijian Lu, and Ling Shao. Polarmix: A general data augmentation technique for lidar point clouds. *CoRR*, abs/2208.00223, 2022. 2

[50] Bin Xu and Zhenzhong Chen. Multi-level fusion based 3d object detection from monocular images. In *CVPR*, pages 2345–2353, 2018. 1

[51] Yan Yan, Yuxing Mao, and Bo Li. Second: Sparsely embedded convolutional detection. *Sensors*, 18(10):3337, 2018. 2

[52] Fisher Yu, Dequan Wang, Evan Shelhamer, and Trevor Darrell. Deep layer aggregation. In *CVPR*, pages 2403–2412, 2018. 11

[53] Wenwei Zhang, Zhe Wang, and Chen Change Loy. Exploring data augmentation for multi-modality 3d object detection. *arXiv preprint arXiv:2012.12741*, 2020. 2, 7

[54] Yunpeng Zhang, Jiwen Lu, and Jie Zhou. Objects are different: Flexible monocular 3d object detection. In *CVPR*, pages 3289–3298, 2021. 1

[55] Yinmin Zhang, Xinzhu Ma, Shuai Yi, Jun Hou, Zhihui Wang, Wanli Ouyang, and Dan Xu. Learning geometry-guided depth via projective modeling for monocular 3d object detection. *arXiv preprint arXiv:2107.13931*, 2021. 1

[56] Barret Zoph, Ekin D Cubuk, Golnaz Ghiasi, Tsung-Yi Lin, Jonathon Shlens, and Quoc V Le. Learning data augmentation strategies for object detection. In *ECCV*, pages 566–583. Springer, 2020. 2

## A. Preliminaries of NeRF

The 3D scene can be represented with the NeRF model, and the neural network is utilized to map a point position $\boldsymbol{x} \in \mathbb{R}^3$ and a view direction $\boldsymbol{d} \in \mathbb{R}^3$ to the corresponding color $\boldsymbol{c} \in \mathbb{R}^3$ and volume density $\sigma$ [27]. We apply the voxel grid to represent the scene considering the low computational cost of voxel-based NeRF [41]. The density voxel grid $\boldsymbol{V}_{\text{density}}$ and feature voxel grid $\boldsymbol{V}_{\text{color}}$ with a shallow MLP are adopted to represent the scene geometry and appearance, respectively. Given input queries $\boldsymbol{x}$ and $\boldsymbol{d}$, the outputs are obtained with the interpolation

$$\begin{aligned} \sigma &= \text{inter}(\boldsymbol{x}, \boldsymbol{V}_{\text{density}}) \\ \boldsymbol{c} &= \text{MLP}_\theta(\text{inter}(\boldsymbol{x}, \boldsymbol{V}_{\text{color}}), \boldsymbol{x}, \boldsymbol{d}) \end{aligned} \tag{7}$$

To render the image, the pixel color $\boldsymbol{C}(\boldsymbol{r})$ along the camera ray $\boldsymbol{r}(t) = \boldsymbol{r_0} + t\boldsymbol{d}$ is approximated by the volume rendering

$$\boldsymbol{C}(\mathbf{r}) = \int_{t_1}^{t_2} T(t)\sigma(\boldsymbol{r}(t))\mathbf{c}(\boldsymbol{r}(t), \mathbf{d})dt, \tag{8}$$

where $t_1$ and $t_2$ are near and far bounds for sampling points, $\boldsymbol{r_0}$ is the camera origin, and $T(t)$ is accumulated transmittance along the ray from $t_1$ to $t$ defined by

$$T(t) = \exp\left(-\int_{t_1}^{t} \sigma(\boldsymbol{r}(s))ds\right). \tag{9}$$

The NeRF model is trained by minimizing the loss between the rendered pixel color $\boldsymbol{C}(\boldsymbol{r})$ and observed pixel color $\hat{\boldsymbol{C}}(\boldsymbol{r})$ given by

$$\mathcal{L}_{\text{Color}} = \sum_{\boldsymbol{r} \in \mathcal{R}(\mathbf{P})} \|\hat{\boldsymbol{C}}(\boldsymbol{r}) - \boldsymbol{C}(\boldsymbol{r})\|_2^2, \tag{10}$$

where $\mathcal{R}(\mathbf{P})$ is the set of rendered rays in a batch.

## B. Implementation Detail of Detectors

For FCOS3D, we utilize a ResNet-101-DCN[18, 8] as the backbone. The model is trained for 24 epochs using the SGD optimizer with an initial learning rate of 1e-4 and a momentum of 0.9. We set the weight decay to 1e-5, and the max norm of gradient clipping to 35. We also adopt a step decay learning rate scheduler with a $0.1\times$ decrease at epoch 20 and 23, along with 1000 iterations of linear warm-up. For SMOKE, we employ a DLA-34[52] as the backbone. We use the Adam optimizer with an initial learning rate 1e-4, and the remaining settings are the same as FCOS3D. For both detectors, their backbones are initialized with ImageNet pre-trained weights. The batch size for training is set to 16.

| 3DAug | DS | LET-AP | LET-APH | LET-APL |
|---|---|---|---|---|
| - | - | 0.585 | 0.573 | 0.393 |
| ✓ | - | 0.584 | 0.572 | 0.394 |
| ✓ | ✓ | **0.590** | **0.578** | **0.403** |

Table 6. **Ablation Study** of Drive-3DAug for FCOS3D on Waymo validation set. 3DAug means we use DVGO [41] for data augmentation. DS means depth supervision.
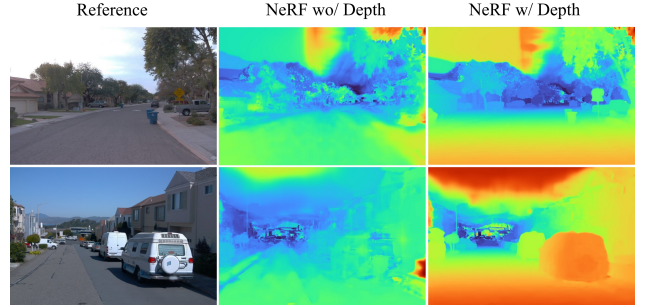


Figure 10. **Visualization of rendered depth map.** The model with depth supervision depicts better performance.

## C. Ablation Study on Depth Supervision.

We qualitatively and quantitatively investigate the effect of depth supervision on background model training and 3D augmentation. Table 6 shows that the 3D augmentation based on background model trained with depth supervision has better performance, with LET-AP (0.590 vs 0.585). Figure 10 shows that NeRF can reconstruct the background with high quality given depth supervision, and the 3D background model quality can be decreased without depth information. Thus, LET-AP, LET-APH and LET-APL on car have a slight decrease with 0.001 for 3D augmentation using background model without supervision.

## D. Visualization of Drive-3DAug

We augment car in Figure 11, indicating that we can generate scene with high quality with Drive-3DAug. Compared with car, pedestrian and cyclist are not rigid body and the size is small, not well applicable for NeRF modelling. We model pedestrian and cyclist as rigid boby in the present study, which can cause the decay of object model performance. As shown in Figure 12, although there exists flaw for augmented pedestrian and cyclist, we can still augment them to improve the detector performance.

## E. Cross-dataset Drive-3DAug

We have reconstructed thousands of background and object models in Waymo and nuScenes dataset. These models can serve as the general model assets, convenient for cre-

ating new driving scenes inside a specific dataset or cross different datasets. As shown in Figure 13, we compose the object models from nuScenes and the background models from Waymo to create new driving scenes. This can further enlarge the diversity of the training data, and we can generate large amounts of data for the study of model generalization across different datasets.

Figure 11. **Visualization of the generated images by Drive-3DAug.** The yellow boxes indicate the newly added cars for the background.



Figure 12. **Visualization of the generated images by Drive-3DAug.** The yellow and red boxes indicate the augmented pedestrian and cyclist, respectively.



Figure 13. **Image Generation cross datasets.** We place the cars from nuScenes on the backgrounds of Waymo. The yellow boxes indicate the augmented cars.