# Improving 3D Object Detection for Pedestrians with Virtual Multi-View Synthesis Orientation Estimation

Jason Ku, Alex D. Pon, Sean Walsh, and Steven L. Waslander

*Abstract*— Accurately estimating the orientation of pedestrians is an important and challenging task for autonomous driving because this information is essential for tracking and predicting pedestrian behavior. This paper presents a flexible Virtual Multi-View Synthesis module that can be adopted into 3D object detection methods to improve orientation estimation. The module uses a multi-step process to acquire the fine-grained semantic information required for accurate orientation estimation. First, the scene's point cloud is densified using a structure preserving depth completion algorithm and each point is colorized using its corresponding RGB pixel. Next, virtual cameras are placed around each object in the densified point cloud to generate novel viewpoints, which preserve the object's appearance. We show that this module greatly improves the orientation estimation on the challenging pedestrian class on the KITTI benchmark. When used with the open-source 3D detector AVOD-FPN, we outperform all other published methods on the pedestrian Orientation, 3D, and Bird's Eye View benchmarks.

## I. INTRODUCTION

Deep neural networks have made remarkable progress on the task of 3D object detection such that they are robust enough to be deployed on autonomous vehicles. The KITTI [1] benchmark showcases the success of 3D object detection methods, especially on the car and cyclist classes, but also highlights areas needing improvement. The benchmark shows that existing 3D detection methods [2], [3] are able to very accurately estimate the orientation of cars and cyclists, with Average Angular Errors (AAE) [4] less than 7° and 20°, respectively, whereas the average error for pedestrians is almost 56°.

In this work, we address the task of 3D pose estimation for pedestrians with a focus on orientation estimation. This task is especially important for autonomous driving as this information is useful for tracking and is vital for predicting pedestrian behavior. Furthermore, it is important to incorporate orientation estimation into the detection pipeline, and not rely on a tracking method for this estimation, as the orientation of pedestrians waiting to cross the street must be inferred without the aid of motion cues.

Image based detection methods have access to rich semantic information from RGB data. To accurately estimate orientation, these methods must extract fine-grained details of objects. However, extracting semantic information is

Jason Ku, Alex D. Pon, Sean Walsh, and Steven L. Waslander are with the Institute for Aerospace Studies, University of Toronto, 4925 Dufferin St, North York, ON, Canada.
`kujason.ku@mail.utoronto.ca,`
`alex.pon@mail.utoronto.ca,`
`sean.walsh@mail.utoronto.ca,`
`stevenw@utias.utoronto.ca`

Fig. 1: **Virtual Multi-View Synthesis.** The core idea of the method is to generate a set of virtual views for each detected pedestrian, and exploit these views in both the training and inference procedures to produce an accurate orientation estimation.

challenging because of the varying scale and appearance of objects caused by the perspective transformation of the 3D scene into an image. Some methods [5], [6] attempt to resolve this issue by using images at multiple scales to extract features. However, object appearance inconsistencies still exist from ROI cropping as noted in [4]. In contrast, we propose to learn fine-grained information by rendering multiple viewpoints of objects through the placement of virtual cameras within the 3D scene at consistent locations relative to each object. Using these generated viewpoints retains a more consistent object appearance, as shown in Fig 1, which facilitates the learning of object orientation in our neural network.

LiDAR methods can take advantage of accurate depth information to achieve robust localization. For the car and cyclist classes, the ratio of the 3D object length and width can be leveraged to simplify the orientation estimation problem [2]. However, LiDAR based methods also face challenges in extracting fine-grained semantic information of objects. The sparsity of LiDAR data limits the operational range of these methods, especially for smaller objects such as pedestrians. At longer distances, fine-grained information is

lost as the sparsity of the LiDAR data becomes so severe that it can become difficult to distinguish between trees, poles, and pedestrians, leading to false positives. Fig. 2 shows that even at shorter distances of 20 and 30 meters, with a high density 64 beam HDL-64E LiDAR, it is difficult for humans to discern meaningful orientation information due to the sparsity of the data. To solve this sparsity issue, we leverage the task of depth completion to generate dense point clouds, which allows for a one-to-one pixel-to-point incorporation of RGB image data. Also, inspired by F-PointNet's [7] use of a 2D detector for accurate classification, we leverage a 2D detector for false positive suppression.

Moreover, orientation estimation performance is dependent on the amount of available training data. The training set must well represent the entire range of possible orientations, but labelled 3D data is expensive and time consuming to acquire [8]. The KITTI [1] dataset has only 4500 pedestrian training instances with 3D labels, making it difficult to train neural networks that normally require much larger amounts of data to achieve good performance. A common solution for the scarcity of pose data is to leverage CAD models and additional annotations for more training data [9], [10]. We, however, do not require the use of additional data sources or labels. To mitigate the low amount of data, we develop a virtual multi-view rendering pipeline to generate novel realistic data from the image and LiDAR inputs. This generated data is incorporated into our network during both training and inference. Detected objects are re-rendered from a set of canonical camera views so that objects maintain a more consistent appearance as compared to using the 2D crops obtained by the common Region of Interest (ROI) crop-and-resize procedure. At inference time, these additional viewpoints are used to determine a more accurate orientation estimate.

To summarize, our contributions in this work are as follows:

- We propose a flexible module that can be adopted into 3D object detection pipelines to improve orientation estimation.
- We solve the issues of limited pose data and the challenges of obtaining fine-grain details from images and LiDAR data in a novel pipeline that uses virtual cameras to generate viewpoints on a colorized depth completed LiDAR point cloud.
- At the time of submission, when incorporated into the detection pipeline of the open-source AVOD-FPN [2] 3D object detector, our method ranks first compared to all other published methods on the KITTI [1] pedestrian Orientation, 3D, and Bird's Eye View benchmarks.
- Compared to other detection methods with comparable orientation estimation performance, our method runs over 8 times faster

## II. RELATED WORK

Orientation (yaw, pitch, roll) represents one of the components defining the 3D pose of an object. The introduction of 3D pose datasets [1], [11], [12], [13] has spurred a myriad of
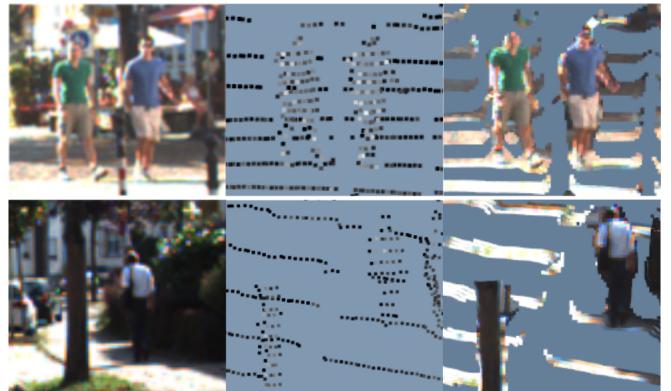


Fig. 2: **Pedestrian Appearances at 20 m (top row) and 30 m (bottom row).** From left to right: RGB image, LiDAR scan colored by intensity, depth completed point cloud colored with corresponding RGB pixels. Even for humans, the classification of objects such as the tree, and the orientations of the pedestrians are not readily apparent in the LiDAR scan. In our method, rich semantic image features are preserved and fused directly with the point cloud, making it much easier to discern this information.

3D pose estimation methods, and as a result many techniques for orientation estimation have been explored.

**Feature Extraction at Multiple Scales.** Previous works have recognized that accurate orientation estimation requires a feature extraction process that captures fine-grained semantic information of objects. Zhang et al. [14] identify that standard Faster R-CNN feature maps are too low resolution for pedestrians and instead use atrous convolutions [15] and pool features from shallower layers. Pyramid structures including image pyramids [16] and feature pyramids [17] have also been leveraged to obtain information from multiple scales. SubCNN [5] use image pyramids to handle scale changes of objects, and [2] highlight the importance of a pyramid structure for small classes such as pedestrians. Moreover, [18], [19], [20] have shown the importance in how methods crop ROI features. Kundu et al. [4] note that the standard ROI crop can warp the appearance of shape and pose, and propose the use of a virtual ROI camera to address this. We instead obtain fine-grained details by using multiple virtual ROI cameras placed canonically around each object's detected centroid. As compared to [4], we not only use 2D RGB data, but also use a 3D point cloud to produce realistic novel viewpoints that maintain consistent object appearance.

**Keypoints and CAD Models.** Using keypoint detections [21], [22], [23], [24] and CAD models [25], [26] have been shown to be effective in gaining semantic understanding of objects of interest. The use of 2D keypoint detections to estimate pose has been well studied as the Perspective-n-Point (PnP) problem with many proposed solutions [27], [28], [29]. More recently, [25], [26] use 3D CAD models and convolutional neural networks (CNNs) to detect keypoints to learn 3D pose. Within the autonomous driving context, Deep MANTA [10] predicts vehicle part coordinates and uses a vehicle CAD model dataset to estimate 3D poses. CAD
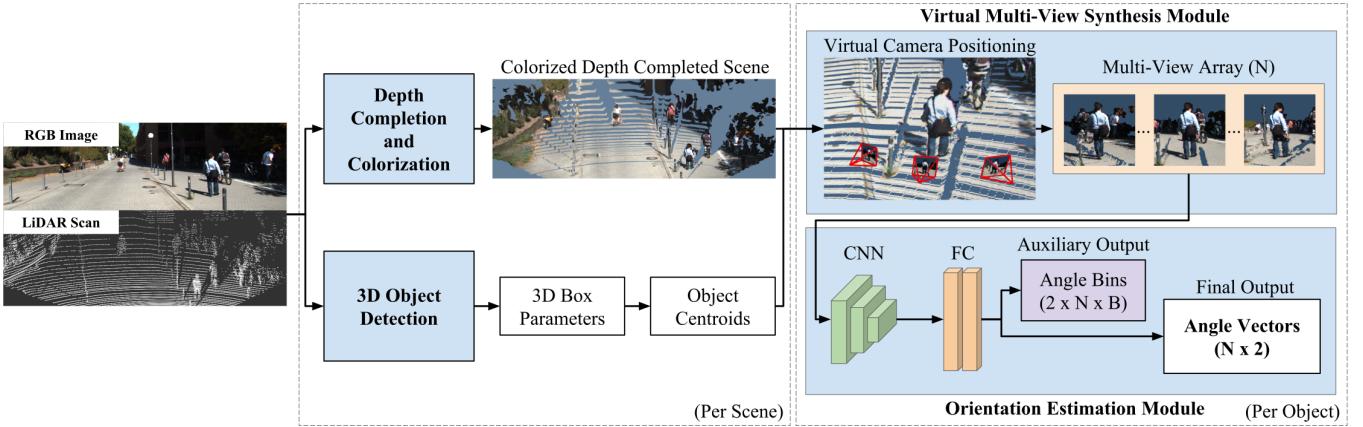
Fig. 3: **Architecture Diagram.** A 3D detector is used to generate 3D detections, which are passed into the Virtual Multi-View Synthesis Module. The module places virtual cameras within the scene represented by a colorized depth completed LiDAR scan to generate $N$ novel viewpoints. Finally, the Orientation Estimation Module predicts the object orientation from the generated views.

models have also been used to create additional ground truth labels. Su et al. [9] argue the scarcity of training data with viewpoint annotations hinders viewpoint estimation performance, so they use 3D models to generate accurate ground truth data. Unlike the above methods, we do not require additional keypoint labels or external CAD datasets. We propose a general pipeline that leverages available data, and with our virtual cameras we generate novel high resolution viewpoints.

**Multi-view Learning.** Using multiple views has previously been shown to be effective in allowing neural networks to learn shape and pose information. Su et al. [30] render multiple views around an object from a CAD dataset, then predict shape based on each view's features. Others [31], [32], [33] use multiple views to ensure a projection consistency to learn shape and pose information. These methods tend to use CAD models, which segment the object of interest from the background, contain full $360°$ shape information, and allows perfect generation of data from any viewpoint. However, our application in autonomous driving only has access to LiDAR scans, which only provides input data from a single direction. We show that we are still able to exploit this data by carefully placing the virtual cameras within a certain operating region to maintain realistically rendered images, as shown in Fig. 4.

**Orientation Representation.** Most similar to our work are 3D pose estimation methods designed for autonomous driving scenarios. These methods have mainly focused on the representation of orientation and designing new loss functions. Pose-RCNN [34] uses a Biternion representation for orientation as recommended by [35]. The monocular 3D object detection method Deep3DBox [36] proposes a *angle bin* formulation that frames orientation estimation as a hybrid classification-regression problem. Here, orientation is discretized into several bins, and the network is tasked to classify the correct bin and to predict a regression offset. This formulation has been adopted by LiDAR methods including [7]. [2] identify an ambiguity issue where identical 3D boxes

are created despite having orientation estimates differing by $\pm\pi$ radians. They solve this issue by parameterizing orientation as an angle vector, while Yan et al. [3] approach this same issue with a sine error loss. We show in our ablation studies (Sec. V-B) that parameterizing orientation as an angle vector while using the discrete-continuous angle bin formulation as an auxiliary loss is most effective.

## III. POSE ESTIMATION FRAMEWORK

Fig. 3 provides an overview of our pipeline for 3D pedestrian pose estimation. Given a set of 3D detections, each parameterized by a centroid $T = (t_x, t_y, t_z)$, dimensions $D = (d_x, d_y, d_z)$, and an Euler angle based orientation $O = (\theta, \phi, \psi)$, the objective is to provide a more accurate estimate of the orientation of the object of interest. For the KITTI benchmark, the pitch and roll are assumed to be zero, while only the yaw, or heading angle, $\theta_{yaw}$ is considered.

The core idea of our method is to generate realistic novel viewpoints of objects from a colored densified point cloud, and to use these viewpoints to extract rich semantic information of objects and thus perform more accurate orientation estimation. To start, we build upon the performance of existing 3D detectors by using 3D detections as centroid proposals. These proposals are processed by the Virtual Multi-View Synthesis module, which renders a set of novel views using a dense point cloud reconstruction. Importantly, these views are created in a set of canonical camera viewpoints, by placing virtual cameras at consistent locations relative to each object's centroid. The object views generated from these virtual cameras better preserves the 3D shape and appearance of the object as compared to taking ROI crops.

For each of the generated novel views, orientation is estimated by passing the views through a CNN followed by an orientation regression output head. The final orientation output is produced by merging the orientation estimates. Lastly, since the 3D detector used in the pipeline is chosen on a basis of having high recall, we use a robust 2D detector for false positive suppression.

## A. Virtual Multi-View Synthesis

**Densified RGB Point Cloud Generation.** We first note that a LiDAR point cloud is simply a sparse representation of the underlying scene, and can be viewed omnidirectionally, with each view providing a unique visual representation of the scene. However the sparsity of LiDAR data only provides a low resolution perspective of the scene. We therefore take the portion of the LiDAR scan corresponding to the visible portion of the scene in the RGB image, and process it through a structure preserving depth completion algorithm [37]. In particular, the LiDAR points are projected into the image using the provided $3 \times 4$ camera projection matrix $P_{cam}$ creating a sparse depth map $D_s$. The depth completion algorithm produces a dense depth map $D_d$, with each pixel representing a depth, that is then re-projected as a 3D point cloud. As the resulting point cloud comes from the dense depth map which has the same resolution as the RGB image, we next infuse the semantic information from the image by coloring each 3D point with its corresponding pixel from the RGB image. This colorized, dense reconstruction of the scene's point cloud resolves the low resolution of the LiDAR scan, and allows for more realistic novel views to be generated that preserve fine-grained semantic information, as shown in Fig. 2.

**Multi-View Generation.** Our goal is to learn fine-grained details of each object. The most straightforward process is to simply take ROI crops of the areas of the image corresponding to each object, but as noted by [4] this common cropping procedure can lead to wildly different appearances of the same object. The use of object-centric coordinate frames has recently been shown to be effective for facilitating learning tasks [7], [38], and as such, we design our module to use canonical viewpoints for estimation.

To create our canonical camera viewpoints, we place virtual cameras in positions equidistant from the object centroid along $N$ equally spaced angle divisions. The camera locations are placed level with the center of the object of interest and between a range of angles, $\rho \in [-\rho_{max}, \rho_{max}]$ to the left and right of the horizontal viewing angle $\alpha$, defined by the ray from the original camera center to the object centroid, as shown in Fig. 4. The minimum and maximum angle of the viewpoints are chosen such that the generated views do not expose parts of the object that are not visible from the original camera viewpoint, which would make the object look unrealistic. $N$ views of the object are generated in a set of canonical viewpoints, which are evenly spaced along an arc of radius $r$ around the object. Each viewpoint generates an $H \times W$ ROI image that preserves appearance, as shown in Fig. 1. These rendered outputs densely fuse the point cloud and image information into a 3 channel RGB format, which allows the use of mature CNN architectures to be used for orientation estimation.

## B. Orientation Estimation

The rendered ROI images are passed through a CNN to produce the final orientation estimations. Several methods [36], [7] use a discrete-continuous loss in the form of $B$
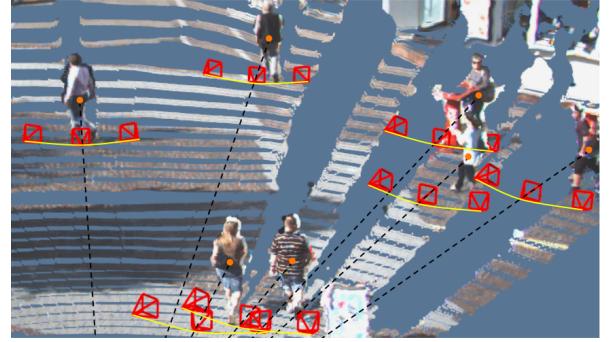


Fig. 4: **Virtual Camera Placement.** Virtual cameras are placed at positions equidistant from each object centroid, with viewpoints ranging from $-25°$ to $25°$ relative to the ray from the original camera center to the object centroid, shown by the dotted black line. Here, only three of the eleven camera positions are shown.

angle bins with regressions within each bin. However, we hypothesize that this divides the training data as each bin will only have a fraction of the total number of training samples to learn from. As shown in Sec. V-B, the network trained in this way does not work as well as when using the angle vector formulation. This is likely due to the small number of pedestrians instances available for training. In our network, the orientation is instead predicted as two values in a vector format, $(x_\theta, y_\theta) = (\cos(\theta), \sin(\theta))$, in the same manner as [2]. This formulation avoids dividing the training samples and handles angle wrapping as each angle $\theta$ is represented as a unique unit vector. At inference time, the ROI images corresponding to each object are processed through the network as a single batch to produce $N$ orientation estimates. The final yaw estimation $\theta_f$ is calculated as the mean of the angle vector yaw estimates as:

$$\theta_f = \tan^{-1} \left( \frac{\frac{1}{N} \sum_{j=1}^{N} y_{\theta_j}}{\frac{1}{N} \sum_{j=1}^{N} x_{\theta_j}} \right). \qquad (1)$$

## C. Final Pose Estimation

The final 3D detection is parameterized by the object centroid $T$, dimensions $D$, and orientation $O$. The centroid and dimension estimates are taken directly from the 3D detections, while the yaw $\theta$ of the orientation is set as the output produced by the Orientation Estimation module.

## D. False Positive Suppression

Pedestrian LiDAR-based 3D detectors' average precision (AP) curves on the KITTI test benchmark show significantly lower precision compared to methods that incorporate image data. This supports the hypothesis that it is significantly harder to discern smaller objects like pedestrians from only LiDAR data, and that false positives are prevalent. To improve detection performance, the 2D boxes from a robust 2D detector are matched with the 2D projections of the 3D detections using an intersection over union (IoU) threshold of 0.4. 3D boxes whose projections do not meet this threshold have their scores reduced, which acts as false

positive suppression. Since the KITTI Average Orientation Similarity (AOS) evaluation is also dependent on 2D detection performance, the projection of each 3D box is replaced with its corresponding 2D box from the 2D detector, which allows AOS to be evaluated.

### E. Training Losses

Multi-task training has shown to be effective in improving performance of many neural networks. Following the discrete-continuous angle bin formulation from [36], we add an auxiliary output layer that produces $B$ heading bins and $B$ angle regressions. The angle bins are trained with a softmax loss while the bin regression and angle vector outputs are trained with a smooth L1 loss. The total loss is computed as:

$$L = \alpha L_{ang} + \beta L_{cls} + \gamma L_{reg}, \quad (2)$$

where $L_{ang}$ is the angle vector loss, $L_{cls}$ is the angle bin classification loss, and $L_{reg}$ is the angle bin regression loss. Losses are weighted such that they converge to values of similar magnitude. We use values of $\alpha = 50$, $\beta = 1$, and $\gamma = 200$. The output from the angle bins are not used in the final estimation as they provide lower performance when used alone. However, incorporating these bins for auxiliary losses in the multi-task training scheme provides additional orientation estimation performance as shown in Sec. V-B.

## IV. IMPLEMENTATION DETAILS

### A. Detector Selection

We use the publicly available AVOD-FPN [2] 3D object detector to generate 3D detections. The selection of this detector is based on the recall analysis in Sec. V-A, which shows that it achieves the highest recall for pedestrians compared to other open-source 3D detectors. The open-source 2D detector MS-CNN [39] is used for false positive suppression. However, our method is compatible with any 3D or 2D detector.

### B. View Synthesis

For each detected object, the camera viewpoints are set to $N = 11$ evenly spaced locations ranging from $\rho \in [-25°, 25°]$ relative to the viewing angle of the object, at a distance $r = 4$ m as shown in Fig. 4. These minimum and maximum angles are empirically chosen to avoid viewing unseen portions of the pedestrians, which are not captured in the LiDAR scan. The ROI images of size $H \times W$ are generated at a resolution of $224 \times 224$. All 11 images are efficiently rendered in parallel using the Visualization Toolkit (VTK).

### C. Orientation Estimation Module

For faster training and convergence, the backbone of the network consists of a modified ResNet-101 [40] architecture, taking only the layers before *conv5_1*, pretrained on the KITTI 2D detection dataset. To reduce GPU memory usage, a 1x1 convolutional layer is applied to the final feature map layer to reduce the depth to 256. The output of this layer

| Method | Recall | | |
|--------|--------|----------|------|
| | Easy | Moderate | Hard |
| SECOND[3] | 86.33 | 80.22 | 73.11 |
| F-PointNet[7] | 95.15 | 88.50 | 80.66 |
| AVOD-FPN[2] | **96.47** | **93.90** | **91.16** |

TABLE I: **3D Pedestrian Detector Recall Analysis.** Recall on moderate pedestrians for top 16 detections at a distance of 0.3 m on the *validation* split. We determine that AVOD-FPN [2] has the highest recall of the available open-source 3D detectors.

| Bins | Vector | OS |
|------|--------|------|
| 4 | - | 0.8817 |
| 8 | - | **0.8940** |
| 12 | - | 0.8892 |
| 16 | - | 0.8625 |
| - | ✓ | 0.9060 |
| 8 | ✓ | **0.9148** |

TABLE II: **Angle Regression Ablation.** Comparison of Orientation Score (OS) performance for different types of angle regression.

| Viewpoint | Views | Time (s) | OS |
|-----------|-------|----------|------|
| ROI Crop | 1 | 0.042 | 0.8909 |
| Virtual Cam | 1 | 0.042 | 0.8692 |
| Virtual Cam | 3 | 0.111 | 0.9001 |
| Virtual Cam | 6 | 0.209 | 0.9074 |
| Virtual Cam | 11 | 0.378 | **0.9148** |
| Virtual Cam | 15 | 0.520 | 0.9061 |

TABLE III: **Virtual Multi-View Ablation.** Effect of ROI versus virtual camera, the number of multi-views on CNN inference time for 16 objects, and orientation score (OS) when used for training and inference.

is flattened and passed through two fully connected layers of size $(512, 256)$, followed by the orientation estimation output heads. The network is trained on batches of 32 ROI images, while inference can run on up to 16 objects (176 ROI images) simultaneously. The network is trained using an Adam optimizer for 20 epochs with an initial learning rate of 0.0001, which is decayed to 0.00005 at 10 epochs.

## V. EXPERIMENTS AND RESULTS

We validate our method on the challenging KITTI Object Detection benchmark, where we evaluate our results on the Orientation, 3D, and Bird's Eye View (BEV) object detection tasks. We compare with current state-of-the-art methods across all tasks, and show performance improvement across all categories. We follow the common training and validation split used by other methods [2], [7], [41], [42], which split the training data at a roughly 1:1 training to validation ratio. As with [2], [7], a separate training split is used for test submission for better generalization. For training, ROI images are generated from the ground truth centroid annotations, and are augmented with contrast, brightness, and left-right flipping. To simulate 3D localization errors from the 3D detection network, 2D translation is sampled from a uniform distribution up to 10 pixels in any direction. In our ablation studies, we validate our method using renders centered around the top 16 scoring pedestrian centroids per scene

| Method | Input | Output | 2D AP | | | AOS | | | OS | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Easy | Moderate | Hard | Easy | Moderate | Hard | Easy | Moderate | Hard |
| 3DOP[41] | RGB+Stereo | 2D | 81.78 | 67.47 | 64.70 | 72.94 | 59.80 | 57.03 | 0.8919 | 0.8863 | 0.8815 |
| SubCNN[5] | RGB | 2D | **83.28** | **71.33** | 66.36 | 78.45 | 66.28 | 61.36 | 0.9420 | 0.9291 | 0.9247 |
| Pose-RCNN[34] | RGB+LiDAR | 2D | 77.53 | 63.40 | 57.49 | 73.95 | 59.90 | 54.27 | 0.9538 | 0.9448 | 0.9440 |
| SECOND[3] | LiDAR | 3D | 65.73 | 55.74 | 49.08 | 51.56 | 43.51 | 38.78 | 0.7844 | 0.7806 | 0.7901 |
| AVOD-FPN[2] | RGB+LiDAR | 3D | 67.32 | 58.42 | 57.44 | 53.36 | 44.92 | 43.77 | 0.7926 | 0.7689 | 0.7620 |
| Ours | RGB+LiDAR | 3D | 81.11 | 70.89 | **67.23** | **78.57** | **67.66** | **63.83** | **0.9687** | **0.9544** | **0.9494** |

TABLE IV: **Pedestrian Pose Estimation.** $AP_{2D}$, $AOS$, and Orientation Similarity ($OS$) on the KITTI *test* set. $OS$ is $AOS/AP_{2D}$. Our pose estimation outperforms all other methods.

| Method | BEV AP | | | 3D AP | | |
|---|---|---|---|---|---|---|
| | Easy | Moderate | Hard | Easy | Moderate | Hard |
| SECOND[3] | 55.10 | 46.27 | 44.76 | 51.07 | 42.56 | 37.29 |
| F-PointNet[7] | 58.09 | 50.22 | 47.20 | 51.21 | 44.89 | 40.23 |
| AVOD-FPN[2] | 58.75 | 51.05 | 47.54 | 50.80 | 42.81 | 40.88 |
| Ours | **61.46** | **51.73** | **47.69** | **53.98** | **45.01** | **41.72** |

TABLE V: **3D Pedestrian Localization and Detection.** $AP_{BEV}$ and $AP_{3D}$ AP on the KITTI *test* set. Our method outperforms the previous state-of-the-art.

generated by AVOD-FPN [2]. Experiments are compared using Orientation Score [36] ($OS = AOS/AP_{2D}$), which is the AOS normalized by detection AP. For inference on the test data, we also use the 3D detections from AVOD-FPN. We follow the standard KITTI evaluation difficulty categories of easy, moderate, and hard to evaluate our estimations.

### A. Object Recall Analysis

For our method, the selection of the 3D detector is based on the need for high recall as a robust 2D detector can provide strong false positive rejection. Recall analysis for pedestrians based on 3D box IoU can be affected greatly by the initial pose estimation since the bounding boxes are very small. Since we only care about the centroid location and we note that errors along the vertical direction are usually minimal, or less than 8 cm on average, we use an alternative recall metric that evaluates the centroid distance at a threshold of 0.3 m along the horizontal and depth directions, $x$ and $z$, to evaluate whether an object has been detected.

In Tab. I, we evaluate the recall of three pre-trained open-source 3D detection methods using the top 16 scoring pedestrian detections per scene on the validation split. We find that the LiDAR only method, SECOND, misses a larger portion of instances, with a recall of 80.22% on the moderate pedestrian instances, while the fusion based AVOD-FPN detector provides the highest recall of 93.90% on moderate pedestrians. F-PointNet [7] also provides a similarly high recall of 88.50% on moderate, but the 2D detector used in their method is not publicly available.

### B. Angle Regression

Tab. II investigates the effect of using each type of orientation representation. The table shows that the angle bin estimation provides better results at 8 bins compared to 12 or 16, which aligns with the intuition that splitting the training data for each bin results in worse performance. However, using only 4 bins requires the network to learn larger regression values, which does not work as well. The angle vector formulation outperforms the angle bin method,

and the final row of the table shows that using the multi-task loss formulation with the angle bins as an auxiliary loss provides the best result.

### C. Effect of Rendered Views

Tab. III shows the influence of the rendered views on the CNN inference time and orientation estimation performance. The basic ROI crop provides better results than a single virtual camera, likely due to the fact that the novel views are being rendered from sparser data. However, with only 3 views, the multi-view formulation outperforms the ROI crop. Performance increases with the number of rendered views, but the improvement comes at the cost of increased inference time on the GPU. However, at 15 views performance no longer improves likely due to the fact the novel views are rendered too close together resulting in repeated data. The final version of the pipeline uses 11 views, which provides the best OS performance. The rendering of 11 views only takes an additional 5 ms over a single view, but each additional view takes 34 ms to inference in the network.

### D. Comparison with State-of-the-Art

We apply our method on the unseen data of the test samples, and submit the results to the online server for evaluation. In Tab. IV, we compare our method with existing methods on the tasks of 2D detection and orientation estimation. On the easy and moderate difficulties, our proposed method achieves comparable 2D AP performance to other 2D methods that also provide orientation estimation, while performing better on hard instances due to the high recall of our selected 3D detector. On the task of orientation estimation, evaluated as AOS on the test server, our method outperforms all other methods on the benchmark across all difficulties, and produces the highest OS. We also compare our results to other 3D detection methods that provide orientation estimates and show significant improvements in orientation estimation performance.

In Tab. V, we compare our results with the existing state-of-the-art 3D detections methods on the tasks of 3D localization $AP_{BEV}$ and 3D detection $AP_{3D}$. As expected, our
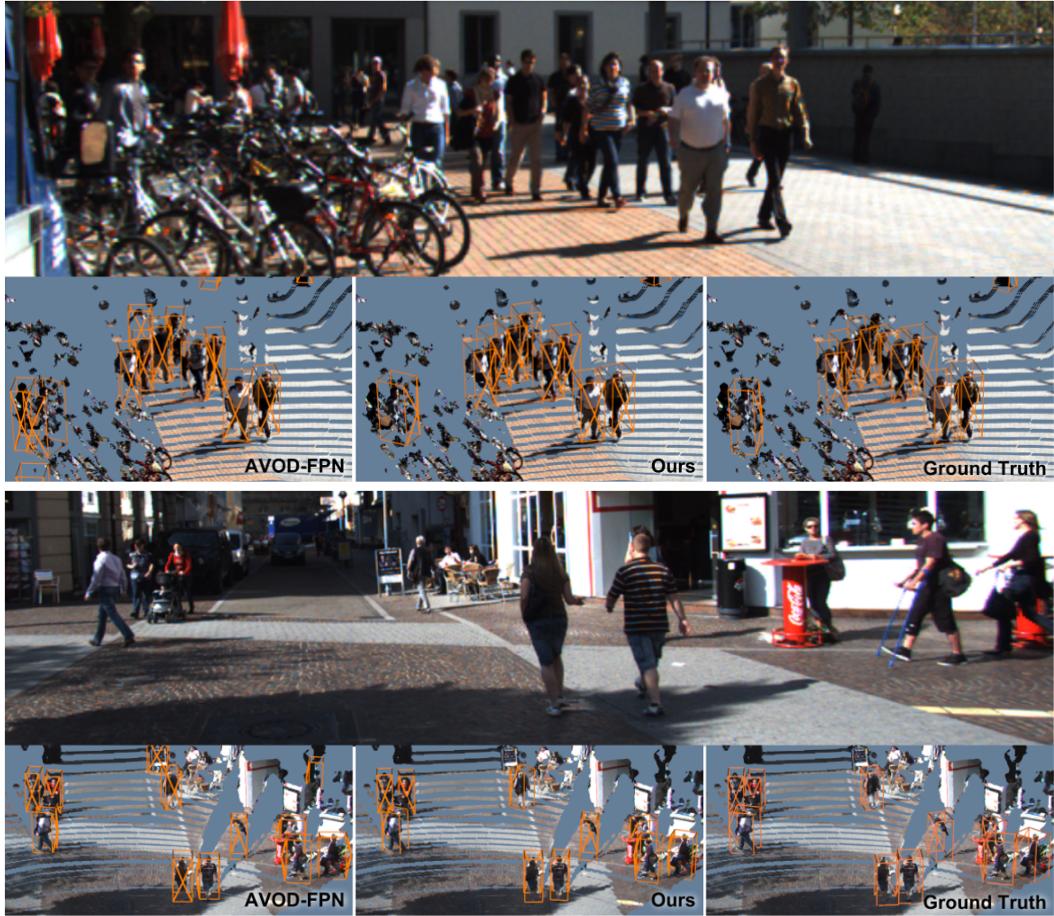
Fig. 5: **Qualitative Results.** We show a comparison of AVOD-FPN detections with and without our orientation module. From left to right: AVOD-FPN [2], Ours, Ground Truth. AVOD-FPN detects all pedestrian instances, but the orientation estimation is poor for several objects and also includes false positives in the detections. Our method estimates orientations that more closely match the ground truth, while also removing false positives.

performance improves upon the base AVOD-FPN detector due to false positive rejection from MS-CNN, as well as the more accurate orientation estimation which increases 3D and BEV IoU due to better 3D box alignment.

### E. Timing

We assess the run time of each component in our proposed method. For depth completion, we use the multi-scale version of [37] which takes 33 ms per LiDAR scan. The View Synthesis module takes 40 ms to render 11 views for 16 objects with VTK. On the KITTI dataset, the CNN inference takes 72 ms per frame on average on a Titan Xp GPU using the AVOD-FPN 3D detections, which equates to 24 ms per object. The AVOD-FPN 3D detector runs at 100 ms per frame. In total, the average time per frame is 245 ms. However, it is important to note that our pipeline does not affect the initial object detection time, which depends on the 3D object detector used. Even so, the full estimation is still 8x faster than other methods with comparable pose estimation performance; SubCNN [5] and Pose-RCNN [34] take 2 s, and 3DOP [41] takes 3 s per frame.

## VI. QUALITATIVE RESULTS

Fig. 5 shows a qualitative comparison of pose estimation results on the KITTI dataset. AVOD-FPN detects all pedestrian instances within the scene. However, orientations are poorly estimated for several instances and there are also false positives. Our module successfully recovers more accurate orientation estimations for the pedestrians within the scene, and suppresses several false positives.

## VII. CONCLUSION

This work presents a method to generate virtual viewpoints on a depth completed point cloud fused with RGB data that addresses the issues of learning semantic information from sparse LiDAR data, and learning from small amounts of pose training data. The generation of virtual views in canonical viewpoints also preserves object appearance. The accurate estimation of pedestrian pose is especially important for safe autonomous driving as the behavior of pedestrians can be more easily predicted. Our module is designed to be incorporated into existing 3D detectors to provide accurate pose estimation, and sets new state-of-the-art results on the KITTI object detection benchmark.

REFERENCES

[1] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *CVPR*, 2012.

[2] J. Ku, M. Mozifian, J. Lee, A. Harakeh, and S. Waslander, "Joint 3d proposal generation and object detection from view aggregation," *IROS*, 2018.

[3] Y. Yan, Y. Mao, and B. Li, "Second: Sparsely embedded convolutional detection," *Sensors*, vol. 18, no. 10, p. 3337, 2018.

[4] A. Kundu, Y. Li, and J. M. Rehg, "3d-rcnn: Instance-level 3d object reconstruction via render-and-compare," in *CVPR*, June 2018.

[5] Y. Xiang, W. Choi, Y. Lin, and S. Savarese, "Subcategory-aware convolutional neural networks for object proposals and detection," in *Applications of Computer Vision (WACV), 2017 IEEE Winter Conference on*. IEEE, 2017, pp. 924–933.

[6] S. Tulsiani and J. Malik, "Viewpoints and keypoints," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1510–1519.

[7] C. R. Qi, W. Liu, C. Wu, H. Su, and L. J. Guibas, "Frustum pointnets for 3d object detection from rgb-d data," in *CVPR*, June 2018.

[8] J. Lee, S. Walsh, A. Harakeh, and S. L. Waslander, "Leveraging pre-trained 3d object detection models for fast ground truth generation," in *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2018, pp. 2504–2510.

[9] H. Su, C. R. Qi, Y. Li, and L. J. Guibas, "Render for cnn: Viewpoint estimation in images using cnns trained with rendered 3d model views," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 2686–2694.

[10] F. Chabot, M. Chaouch, J. Rabarisoa, C. Teulire, and T. Chateau, "Deep manta: A coarse-to-fine many-task network for joint 2d and 3d vehicle analysis from monocular image," in *CVPR*, 2017.

[11] Y. Xiang, W. Kim, W. Chen, J. Ji, C. Choy, H. Su, R. Mottaghi, L. Guibas, and S. Savarese, "Objectnet3d: A large scale database for 3d object recognition," in *European Conference on Computer Vision*. Springer, 2016, pp. 160–176.

[12] Y. Xiang, R. Mottaghi, and S. Savarese, "Beyond pascal: A benchmark for 3d object detection in the wild," in *IEEE Winter Conference on Applications of Computer Vision*. IEEE, 2014, pp. 75–82.

[13] K. Matzen and N. Snavely, "Nyc3dcars: A dataset of 3d vehicles in geographic context," in *Proceedings of the IEEE International Conference on Computer Vision*, 2013, pp. 761–768.

[14] L. Zhang, L. Lin, X. Liang, and K. He, "Is faster r-cnn doing well for pedestrian detection?" in *European conference on computer vision*. Springer, 2016, pp. 443–457.

[15] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2018.

[16] E. H. Adelson, C. H. Anderson, J. R. Bergen, P. J. Burt, and J. M. Ogden, "Pyramid methods in image processing," *RCA engineer*, vol. 29, no. 6, pp. 33–41, 1984.

[17] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2117–2125.

[18] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.

[19] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in Neural Information Processing Systems 28*, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, Eds. Curran Associates, Inc., 2015, pp. 91–99.

[20] B. Jiang, R. Luo, J. Mao, T. Xiao, and Y. Jiang, "Acquisition of localization confidence for accurate object detection," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 784–799.

[21] A. Toshev and C. Szegedy, "Deeppose: Human pose estimation via deep neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 1653–1660.

[22] J. Tompson, R. Goroshin, A. Jain, Y. LeCun, and C. Bregler, "Efficient object localization using convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 648–656.

[23] J. Carreira, P. Agrawal, K. Fragkiadaki, and J. Malik, "Human pose estimation with iterative error feedback," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4733–4742.

[24] K. J. Shih, A. Mallya, S. Singh, and D. Hoiem, "Part localization using multi-proposal consensus for fine-grained categorization," *BMVC*, 2015.

[25] G. Pavlakos, X. Zhou, A. Chan, K. G. Derpanis, and K. Daniilidis, "6-dof object pose from semantic keypoints," in *Robotics and Automation (ICRA), 2017 IEEE International Conference on*. IEEE, 2017, pp. 2011–2018.

[26] J. Wu, T. Xue, J. J. Lim, Y. Tian, J. B. Tenenbaum, A. Torralba, and W. T. Freeman, "Single image 3d interpreter network," in *European Conference on Computer Vision*. Springer, 2016, pp. 365–382.

[27] V. Lepetit, F. Moreno-Noguer, and P. Fua, "Epnp: An accurate o (n) solution to the pnp problem," *International journal of computer vision*, vol. 81, no. 2, p. 155, 2009.

[28] C.-P. Lu, G. D. Hager, and E. Mjolsness, "Fast and globally convergent pose estimation from video images," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 22, no. 6, pp. 610–622, 2000.

[29] A. Ansar and K. Daniilidis, "Linear pose estimation from points or lines," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 25, no. 5, pp. 578–589, 2003.

[30] H. Su, S. Maji, E. Kalogerakis, and E. Learned-Miller, "Multi-view convolutional neural networks for 3d shape recognition," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 945–953.

[31] S. Tulsiani, T. Zhou, A. A. Efros, and J. Malik, "Multi-view supervision for single-view reconstruction via differentiable ray consistency," in *CVPR*, vol. 1, no. 2, 2017, p. 3.

[32] S. Tulsiani, A. A. Efros, and J. Malik, "Multi-view consistency as supervisory signal for learning shape and pose prediction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2897–2905.

[33] P. Mandikal, M. Agarwal, R. V. Babu *et al.*, "Capnet: Continuous approximation projection for 3d point cloud reconstruction using 2d supervision," *arXiv preprint arXiv:1811.11731*, 2018.

[34] M. Braun, Q. Rao, Y. Wang, and F. Flohr, "Pose-rcnn: Joint object detection and pose estimation using 3d object proposals," in *2016 IEEE 19th International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2016, pp. 1546–1551.

[35] L. Beyer, A. Hermans, and B. Leibe, "Biternion nets: Continuous head pose regression from discrete training labels," in *German Conference on Pattern Recognition*. Springer, 2015, pp. 157–168.

[36] A. Mousavian, D. Anguelov, J. J. Flynn, and J. Kosecka, "3d bounding box estimation using deep learning and geometry," *CVPR*, pp. 5632–5640, 2017.

[37] J. Ku, A. Harakeh, and S. L. Waslander, "In defense of classical image processing: Fast depth completion on the cpu," in *CRV*, 2018.

[38] K. Rematas, I. Kemelmacher-Shlizerman, B. Curless, and S. Seitz, "Soccer on your tabletop," in *CVPR*, 2018.

[39] Z. Cai, Q. Fan, R. Feris, and N. Vasconcelos, "A unified multi-scale deep convolutional neural network for fast object detection," in *ECCV*, 2016.

[40] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *CVPR*, pp. 770–778, 2016.

[41] X. Chen, K. Kundu, Y. Zhu, A. Berneshawi, H. Ma, S. Fidler, and R. Urtasun, "3d object proposals for accurate object class detection," in *NIPS*, 2015.

[42] J. Ku, A. D. Pon, and S. L. Waslander, "Monocular 3d object detection leveraging accurate proposals and shape reconstruction," in *CVPR*, 2019.