

# Monocular 3D Object Detection via Geometric Reasoning on Keypoints

Ivan Barabanau  
Skoltech

ivan.barabanau@skoltech.ru

Alexey Artemov  
Skoltech

a.artemov@skoltech.ru

Evgeny Burnaev  
Skoltech

e.burnaev@skoltech.ru

Vyacheslav Murashkin  
Yandex

vmurashkin@yandex-team.ru

## Abstract

*Monocular 3D object detection is well-known to be a challenging vision task due to the loss of depth information; attempts to recover depth using separate image-only approaches lead to unstable and noisy depth estimates, harming 3D detections. In this paper, we propose a novel keypoint-based approach for 3D object detection and localization from a single RGB image. We build our multi-branch model around 2D keypoint detection in images and complement it with a conceptually simple geometric reasoning method. Our network performs in an end-to-end manner, simultaneously and interdependently estimating 2D characteristics, such as 2D bounding boxes, keypoints, and orientation, along with full 3D pose in the scene. We fuse the outputs of distinct branches, applying a reprojection consistency loss during training. The experimental evaluation on the challenging KITTI dataset benchmark demonstrates that our network achieves state-of-the-art results among other monocular 3D detectors.*

## 1. Introduction

The success of autonomous robotics systems, such as self-driving cars, largely relies on their ability to operate in complex dynamic environments; as an essential requirement, autonomous systems must reliably identify and localize non-stationary and interacting objects, *e.g.* vehicles, obstacles, or humans. In its simplest formulation, localization is understood as an ability to detect and frame objects of interest in 3D bounding boxes, providing their 3D locations in the surrounding space. Crucial to the decision-making process is the accuracy of depth estimates of the 3D detections.

Depth estimation could be approached from both hardware and algorithmic perspectives. On the sensors end, laser scanners such as LiDAR devices have been extensively used to acquire depth measurements sufficient for 3D detection in many cases [49, 3, 16, 51, 15, 48]. However, point clouds produced by these expensive sensors are sparse, noisy and massively increase memory footprints with millions of 3D points acquired per second. In contrast, image-based 3D detection methods offer savings on CPU and memory consumption, use cheap onboard cameras, and work with a wealth of established detection architectures (*e.g.*, [24, 35, 36, 19, 11, 20, 23]), yet they require sophisticated algorithms for depth estimation, as raw depth cannot be accessed anymore.

Recent research on monocular 3D object detection relies on separate dense depth estimation models [33, 46], but depth recovery from monocular images is naturally ill-posed, leading to unstable and noisy estimates. In addition, in many practical instances, *e.g.*, with sufficient target resolution or visibility, dense depth estimation might be redundant in context of 3D detection. Instead, one may focus on obtaining sparse but salient features, such as 2D keypoints, that are well-known visual cues often serving as geometric constrains in various vision tasks such as human pose estimation [34, 26, 27, 28] and more general object interpretation [13, 43].

Motivated by this observation, in this paper we propose a novel keypoint-based approach for 3D object detection and localization from a single RGB image. We build our model around 2D keypoint detection in images and complement it with a conceptually simple geometric reasoning framework, establishing correspondences between the detected 2D keypoints and their 3D counterparts defined on surfaces

of 3D CAD models. The framework operates under the general assumptions, assuming the camera intrinsic parameters are given, and retrieves depth of closest keypoint instance, thereby "lifting" 2D keypoints to 3D space; the remaining 3D keypoints and the final 3D detection are assembled in a similar way. Our approach does not require keypoint-annotated labeled images, but instead relies on a multi-task reprojection consistency loss, allowing for robust 2D keypoint detection. Thus, our model is end-to-end trainable.

In summary, our contributions are as follows:

- We propose a novel deep learning-based framework for monocular 3D object detection, combining well-established region-based detectors and a geometric reasoning step over keypoints.
- We describe an end-to-end training scheme for this framework, using a dataset of real-world images and a collection of 3D CAD models, annotated with 3D keypoints.

The rest of this paper is organized as follows. In Section 2, we review the related work on object detection, mostly in the context of self-driving and robotics applications. In Section 3, we describe our proposed monocular 3D object detection approach, and in Section 4, its experimental evaluation using the standard KITTI benchmark. We conclude in Section 5 with a discussion of our results.

## 2. Related work

**2D object detection.** 2D object detection is an extensively studied vision task, with a body of research devoted to both algorithms [35, 24, 20, 36, 11, 23] and benchmarks [5, 6, 21, 8, 7]. Traditionally, object detectors operate in two stages, with the first stage selecting object candidates [41, 54, 36] and the second stage operating as a discriminator and refinement model, rejecting bad proposals [36, 11, 4]. Due to the introduction of novel backbone architectures [45] and losses [20], such approaches have attained top results in a number of benchmarks. In the context of robotics, single-stage detectors such as YOLO [35], SSD [24] and RetinaNet [20] are of particular interest, however, they offer inferior performance and are not straightforward to extend to related tasks such as instance segmentation [11, 23] and keypoint detection [11].

**3D object detection.** Recently, novel deep learning architectures operating directly on unstructured point clouds have been proposed [31, 32, 42, 10, 14, 17], offering

the possibility to develop corresponding 3D object detectors [30, 51, 47]. However, such approaches require expensive sensing equipment (LiDARs) and commonly process point cloud data coupled with RGB data. Some depth-based approaches operate over voxel-grid representations of the point clouds, leveraging the existing convolutional architectures [49, 16, 51, 48], while other methods fuse depth features with birds-eye-view (BEV) and image features [3, 15, 30].

**Monocular 3D object detection.** The most relevant to our work is research on monocular 3D object detection, that is well-known to be a challenging vision task. Deep3DBox [29] relies on a set of geometric constraints between 2D and predicted 3D bounding boxes and reduces 3D object localization problem to a linear system of equations, fitting 3D box projections into 2D detections. Their approach relies on a separate linear solver; in contrast, our model is end-to-end trainable and does not require external optimization. Mono3D [2] extensively samples 14K 3D bounding box proposals per image and evaluates each, exploiting semantic and image-based features. In contrast, our approach does not rely on an exhaustive sampling in 3D space, bypassing a significant computational overhead. OFT-Net [37] introduces an orthographic feature transform which maps RGB image features into a birds-eye-view representation through a 3D scene grid, solving the perspective projection problem. However, back-projecting image features onto 3D grid results in a coarse feature assignment. Our approach detects 2D keypoints with sufficient precision, avoiding any additional discretization. MonoGRNet [33] directly deals with depth estimation from a single image, training an additional sub-network to predict the  $z$ -coordinate of each 3D bounding box. [46] exploit a similar approach, estimating disparity using a stand-alone pre-trained MonoDepth network [9]. Both these methods rely on the non-trainable depth estimation networks, which introduce a computational overhead; in contrast, our approach jointly estimates object 2D bounding-box and 3D pose in a fully trainable manner, not requiring a dense depth prediction.

Perhaps, the most similar approach to ours is [1], which utilizes 3D CAD models, along with predicting 2D keypoints. However, their network only models 2D geometric properties and aims at matching the predictions to one of the CAD shapes, while 3D pose estimation is postponed for the inference step. They additionally exploit extensive annotations of keypoints in their 3D models. In contrast, we

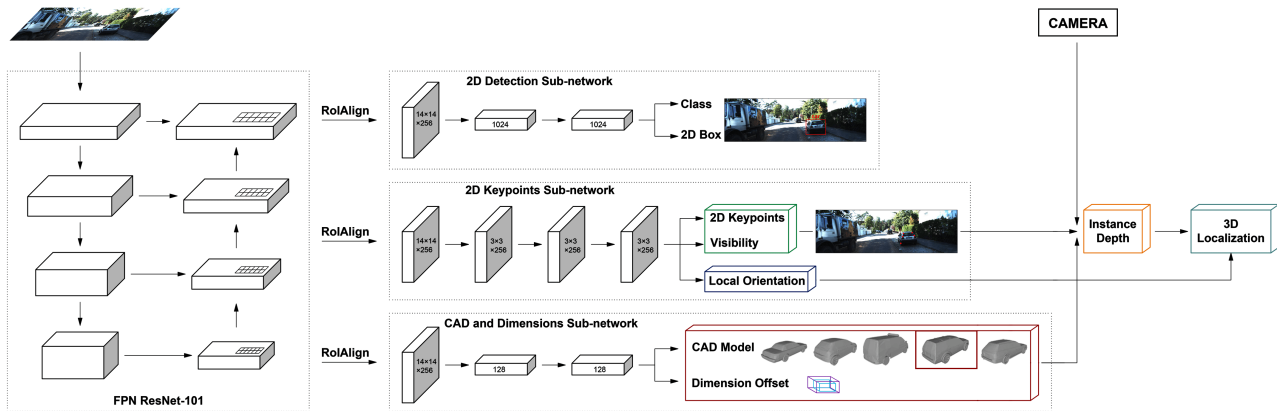


Figure 1: An overview of our monocular 3D detection architecture. We start with a universal backbone network (see use Mask R-CNN [11]) and complement it with three sub-networks: 2D object detection sub-network, 2D keypoints regression sub-network, and dimension regression sub-network. The network is trained end-to-end using a multi-task loss function.

only annotate 14 keypoints per each of the five 3D models and exploit them in a geometric reasoning module to bridge the gap between 2D and 3D worlds, which allows us to deal with 3D characteristics during training in an end-to-end manner.

### Keypoints estimation and 3D representations.

Keypoint-based representations are a common mechanism of encoding 3D geometric structure in objects, that have proven themselves as powerful visual cues for tasks such as pose estimation [34, 26, 27, 28], fine pose prediction [18], 3D reconstruction [38], shape alignment [18, 25], to name a few. A commonly used approach is to learn a set of keypoint detectors, followed by some post-processing to assemble their predictions into a geometric model. However, obtaining sufficient ground-truth for training keypoint detectors is a challenging task. One may manually annotate 3D keypoints of objects in real images, but this is labor-intensive and often inaccurate. Other directions involve active shape modeling [34, 22] and shape alignment with wireframe [52, 53, 50] and 3D CAD models [1]. For human pose estimation, another option could be motion capture of joint locations [26] Recently, latent modeling approaches have been proposed to learn optimal sets of keypoints without direct supervision [43, 39]. Our keypoint detection approach bears similarity to [1] as we utilize 3D CAD models and align them to sensor measurements, but offers labor savings since we only annotate 14 keypoints per each of the

five CAD models.

## 3. 3D Object Detection Framework

Given a single RGB image, our goal is to localize target objects in the 3D scene. To do this, we propose an end-to-end trainable CNN-based framework that accepts a single RGB image as input and outputs a set of 3D detections. Each target object is defined by its class and 3D bounding box, parameterized by 3D center coordinates  $\mathbf{C} = (c_x, c_y, c_z)$  in a camera coordinate system, global orientation  $\mathbf{R} = (\theta, \cdot, \cdot)$ , and dimensions  $\mathbf{D} = (w, h, l)$ , standing for width, height and length, respectively (we don't correct for truncation or occlusion when defining object sizes). We parameterize global object orientation with the yaw angle  $\theta$  only, which is a commonly adopted premise when dealing with objects in the road scenes [29, 30].

Our proposed framework comprises two sub-modules, each of which operates on the characteristics living in either 2D (image) or 3D (world) space. From the 2D perspective, each object of interest, cropped by its predicted 2D bounding box, is provided with 2D keypoints and their respective visibility states. On the 3D side, object dimensions, 3D CAD model, and local orientation are predicted. The gap between the two spaces is bridged by the geometric reasoning module computing instance depth, global orientation, and the final 3D detection.

Our implementation takes advantage of the generality of the state-of-the-art Mask R-CNN architecture [11], view-

ing it as a universal backbone network extensible to adjacent problems, and complement it with three sub-networks: 2D object detection sub-network, 2D keypoints regression sub-network, and dimension regression sub-network. The whole system represents an end-to-end trainable network, depicted in Figure 1, with sub-networks initially trained independently, switching further to joint training via introduced multi-task consistency reprojection loss function on the projected 3D keypoints and 3D bounding box corners.

**2D object detection.** For 2D detection, we follow the original Mask R-CNN architecture [11], which includes Feature Pyramid Network (FPN) [19], Region Proposal Network (RPN) [36] and RoIAlign module [11]. The RPN generates 2D anchor-boxes with a set of fixed aspect ratios and scales throughout the area of the provided feature maps, which are scored for the presence of the object of interest and adjusted. The spatial diversity of the proposed locations is processed by the RoIAlign block, converting each feature map, framed by the region of interest, into a fixed-size grid, preserving an accurate spatial location through bilinear interpolation. Followed by fully connected layers, the network splits into two feature sharing branches for the bounding box regression and object classification. During training, we utilize smooth L1 and cross-entropy loss for each task respectively, as proposed by [36]. Though we do not directly utilize the predicted 2D bounding boxes, we have experimentally observed the 2D detection sub-network to stabilize training.

**2D keypoint detection.** We predict coordinates and a visibility state for each of the manually-chosen 14 keypoints  $\mathbf{K} = \{(x_i, y_i)_{i=1}^{14}\}$  (c.f. Figure 3 for details on our choice of 3D keypoints). Unlike the parameterization suggested in [11, 40], we directly regress on 2D coordinates of keypoints. The visibility state, determined by the occlusion and truncation of an instance, is a binary variable, and no difference between occluded, self-occluded and truncated states is made. Adding visibility estimation helps propagate information during training for visible keypoints only and acts as an auxiliary supervision for orientation sub-network. During training, similar to our 2D object detection sub-network, we minimize the multi-task loss combining smooth L1 loss for coordinates regression and cross-entropy loss for visibility state classification, defined

as:

$$\begin{aligned} \mathcal{L}_{\text{coord}} &= \sum_{k=1}^K \mathbb{1}_k \cdot L_1^{\text{smooth}}(k_{(x,y)}^{\text{gt}}, k_{(x,y)}^{\text{pred}}) \\ \mathcal{L}_{\text{vis}} &= - \sum_{k=1}^K \sum_{j=1}^2 v_{kj}^{\text{gt}} \log v_{kj}^{\text{pred}} \\ \mathcal{L}_{\text{kp}} &= \mathcal{L}_{\text{coord}} + \mathcal{L}_{\text{vis}}, \end{aligned} \quad (1)$$

where  $\mathbb{1}_k$  is the visibility indicator of  $k$ -th keypoint, while  $k_{(x,y)}^{\text{gt}}$  and  $k_{(x,y)}^{\text{pred}}$  denote ground-truth and predicted 2D coordinates, normalized and defined in a reference frame of a specific feature map after RoI-alignment. Similarly,  $v_{kj}^{\text{gt}}$  is the ground truth visibility status, while  $v_{kj}^{\text{pred}}$  is the estimated probability that keypoint  $k$  is visible.

### 3D dimension estimation and geometric classification.

To each annotated 3D instance in the dataset, we have assigned a 3D CAD model out of a predefined set of 5 templates, obtaining 5 distinct geometric classes of instances. Used templates are presented on Figure 2. The assignment has been made based on the width, length and height ratios only. For each geometric class, we have computed mean dimensions  $(\mu_w, \mu_h, \mu_l)$  over all assigned annotated 3D instances.



Figure 2: The 5 geometric classes of instances in our work are represented by 5 3D CAD models with strongly distinct aspect ratios.

During training the 3D dimension estimation and geometric class selection sub-network, we utilize a multi-task loss combining cross-entropy loss (for the geometric class selection) and a smooth  $L_1$  loss for dimension regression. Instead of regressing the absolute dimensions, we predict the differences  $\mathbf{D}_{\text{offset}} = (\Delta w, \Delta h, \Delta l) = (w - \mu_w, h - \mu_h, l - \mu_l)$  from the mean dimensions in the log-space:

$$\mathcal{L}_d(\mathbf{D}^{\text{gt}}, \mathbf{D}^{\text{pred}}) = L_1^{\text{smooth}}(\log(\mathbf{D}_{\text{offset}}^{\text{gt}} - \mathbf{D}_{\text{offset}}^{\text{pred}})) \quad (2)$$

where  $\mathbf{D}_{\text{offset}}^{\text{gt}}$  and  $\mathbf{D}_{\text{offset}}^{\text{pred}}$  represent the ground truth and predicted offsets to the class mean values along each dimension, respectively.

**Reasoning about instance depth.** We define instance depth as the depth  $Z$  of a vertical plane passing through the

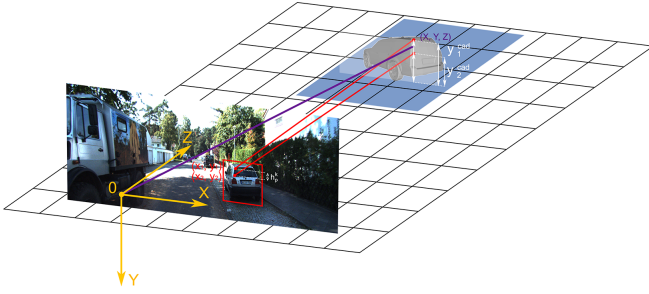


Figure 3: Geometric reasoning about instance depth. We use camera intrinsic parameters, predicted 2D keypoints and dimension predictions to "lift" the keypoints to 3D space.

two closest of visible keypoints, defined in the camera reference frame. To compute this depth value, we use predicted 2D keypoints, instance height (in meters), and its geometric class. First, we select two keypoints  $(x_1, y_1)$  and  $(x_2, y_2)$  in the image and compute their  $y$ -difference  $h_p = |y_1 - y_2|$ . We then select the corresponding two keypoints  $(x_1^{\text{cad}}, y_1^{\text{cad}})$  and  $(x_2^{\text{cad}}, y_2^{\text{cad}})$  in cad model reference frame and compute their height ratio  $r_{\text{cad}} = y_1^{\text{cad}}/y_2^{\text{cad}}$ . Finally, the distance to the object  $Z$  is defined from the pinhole camera model:

$$Z = f \cdot \frac{r_{\text{cad}} \cdot h}{h_p} \quad (3)$$

where  $f$  is a focal length of the camera, known for each frame. Figure 3 illustrates this computation. Depth coordinate  $Z$  allows to retrieve the remaining 3D location coordinates of one of the selected keypoints, using the back-projection mapping:

$$X = Z \cdot \frac{x_{\{1,2\}} - p_x}{f}, \quad Y = Z \cdot \frac{y_{\{1,2\}} - p_y}{f} \quad (4)$$

where  $(p_x, p_y)$  are the camera principal point coordinates in pixels.

**Orientation estimation.** Direct estimation of orientation  $\mathbf{R}$  in a camera reference frame is not feasible, as the region proposal network propagates the context within the crops solely, cutting off the relation of the crop to the image plane. Inspired by [29], we represent the global orientation as a combination of two rotations with azimuths defined as:

$$\theta = \theta_{\text{loc}} + \theta_{\text{ray}} \quad (5)$$

where  $\theta_{\text{loc}}$  is the object's local orientation within the region of interest, and  $\theta_{\text{ray}}$  is a ray direction from the camera to the object center, directly found from the 3D location coordinates. We estimate  $\theta_{\text{loc}}$  using a modification of the Multi-Bin approach [29]. Specifically, instead of splitting the objective into angle confidence and localization parts, we discretize the angle range from  $0^\circ$  to  $360^\circ$  degrees into 72 non-overlapping bins and compute the probability distribution over this set of angles by a softmax layer. We train the local orientation sub-network using cross-entropy as a loss function. To obtain the final prediction for  $\theta_{\text{loc}}$ , we utilize the weighted mean of the bins medians ( $WM(\theta_{\text{loc}})$ ), adopting the softmax output as the weights. Given 3D location coordinates  $(X, Z)$  of one of the keypoints and the weighted mean local orientation  $WM(\theta_{\text{loc}})$ , the global orientation is defined as follows:

$$\theta = WM(\theta_{\text{loc}}) + \arctan\left(\frac{X}{Z}\right). \quad (6)$$

**3D object detection.** To obtain the center  $\mathbf{C}$  of the final 3D bounding box, we use the global orientation  $\mathbf{R}$  and the distance between the keypoint and the object center. For a particular CAD model, given the weight, height and length ratio between the selected keypoint and the object center  $\mathbf{r}_{\text{cad}} = (x_{\text{cad}}^1/x_{\text{cad}}^2, y_{\text{cad}}^1/y_{\text{cad}}^2, z_{\text{cad}}^1/z_{\text{cad}}^2)$  estimated object dimensions  $\mathbf{D}$  and global orientation  $\mathbf{R}$ , the location  $\mathbf{C}$  is predicted as

$$\mathbf{C} = (X, Y, Z) \pm \mathbf{R} \cdot \mathbf{D} \odot \mathbf{r}_{\text{cad}} \quad (7)$$

where  $\odot$  stands for an element-wise product. Depending on the selected keypoint position (left or right, back or front, top or bottom side of the object), a sign is chosen for each dimension.

**Multi-head reprojection consistency loss.** Except for shared convolutional backbone, each sub-network is independent of its neighbors and unaware of other predictions, though the geometric components are strongly interrelated. To provide consistency between the network branches, we introduce a loss function which integrates all the predictions. 3D coordinates in a CAD model coordinate system of the keypoints from the set  $\mathbf{K}$  are scaled using  $\mathbf{D}$ , rotated using  $\mathbf{R}$ , translated using  $\mathbf{C}$ , and back-projected into the image plane via camera projection matrix to obtain 2D keypoint coordinates and compare with the ground truth values. A similar approach is applied to the eight corners of the 3D bounding box obtained from 3D detection and orientation estimates, to ensure that they fit tightly into the ground truth 2D bounding box after back-projection. In all cases, we use

Method	IoU = 0.5			IoU = 0.7		
	Easy	Moderate	Hard	Easy	Moderate	Hard
Mono3D	25.19	18.20	15.52	2.53	2.31	2.31
OFT-Net	-	-	-	4.07	3.27	3.29
MonoGRNet	50.51	<b>36.97</b>	<b>30.82</b>	13.88	<b>10.19</b>	<b>7.62</b>
MF3D	47.88	29.48	<u>26.44</u>	10.53	5.69	<u>5.39</u>
Ours	48.81	30.17	20.07	11.91	6.64	4.28
Ours (+loss)	<b>50.82</b>	<u>31.28</u>	20.21	<b>13.96</b>	<u>7.37</u>	4.54

Table 1: **3D detection performance:** Average Precision of 3D bounding boxes on KITTI val set. The best score is in bold, the second best underlined. Ours (+loss) indicates a base network setup, trained with a consistency reprojection loss.

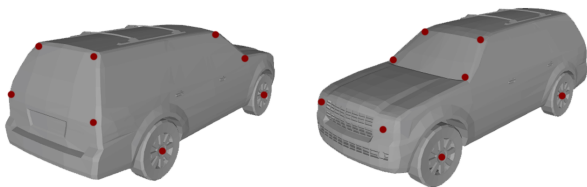


Figure 4: An example of our annotation of 3D keypoints. We label 4 centers of wheels, 8 corners of front and back windshields per each, and 2 centers of headlights. Per each 3D CAD model, we form the annotation so as to ensure that the two keypoints on each of the windshields form near-vertical lines. This is required to ensure accuracy of our geometric reasoning framework.

the smooth L1 loss during training.

## 4. Experiments

### 4.1. Experimental setup

**Dataset.** We train and evaluate our approach using the KITTI 3D object detection benchmark dataset. For the sake of comparison with state-of-the-art methods, we follow the setup presented in [2], which provides 3712 and 3769 images for training and validation, respectively, along with the camera calibration data. To extend KITTI dataset with assignment of geometric classes using CAD models and keypoints 2D coordinates, we employ the approach and data provided in [44]. Depending on the ratios between height, length and width, each car instance is assigned with one out of 5 CAD model classes from a predefined set of CAD templates, presented on Figure 1. We manually annotated each CAD model with the keypoint locations. Figure 3 displays an example of the annotated keypoints, most of which are a common choice [6] due to their interpretability, such as the

car’s edges, carcass, etc.; we also included corners of windshields to deal with the height of each instance. To obtain 2D coordinates of the keypoints, we back-projected CAD models from 3D space to the image plane using ground truth location, dimension and rotation values. Simultaneous projection of all 3D CAD models on a scene provides us with a depth ordering mask, allowing for defining the visibility state of each keypoint.

**Network architecture.** We utilize Mask R-CNN with a Feature Pyramid Network [19], based on a ResNet-101 [12] as our backbone network for the multi-level feature extraction. Instead of the higher resolution  $14 \times 14$  and  $28 \times 28$  feature maps in the original architecture, we stack the same amount of  $3 \times 3$  kernels, followed by a fully connected layer to predict 2D normalized coordinates and visibility states for each of the 14 keypoints. From the same feature maps, we branch a fully-connected layer predicting local orientation in bins of  $5^\circ$  each, totaling 72 output units. The feature sharing between keypoints and local orientation was found crucial for network performance, as both characteristics imply similar geometric reasoning. In parallel to 2D detection and 2D keypoints estimation, we create a sub-network of a similar architecture for dimension regression and classification into geometric classes. The remaining components, including RPN, RoIAlign, bounding box regression and classification heads, are implemented following the original Mask R-CNN design. For instance depth retrieval we use only four pairs of keypoints: corners of the front and rear windows. Other keypoints are used for additional supervision in consistency loss calculation during training.

**Training our model.** We set hyperparameters following Mask R-CNN work [11]. The RPN anchor set covers five scales, adjusting them to the values of 4, 8, 16, 32, 64, and

three default aspect ratios. Each mini-batch consists of 2 images, producing 256 regions of interest, with a positive to negative samples ratio set 1:3, to achieve class sampling balance during training. Any geometric augmentations over the images are omitted, solely applying image padding to meet the network architecture requirements. ResNet-101 is initialized with the weights pre-trained on Imagenet [5], and frozen during further training steps. We first train the 2D detection and classification sub-network for 100K iterations, adopting Adam optimizer with a learning rate of  $10^{-4}$  throughout the training, setting weight decay of 0.001 and momentum of 0.9. Then 2D keypoints and local orientation are trained for 50K iterations. Finally, enabling the multi-head consistency loss, the whole network is trained in an end-to-end fashion for 50K iterations. We combine losses from all of the network outputs, weighting them equally.

**Evaluation metrics.** We evaluate the network under the conventional KITTI benchmark protocol, which enables comparison across approaches. Car category is the sole subject of our focus. By default, KITTI settings require evaluation in 3 regimes: easy, moderate and hard, depending on the instance difficulty of a potential detection. 3D bounding box detection performance implies 3D Average Precision (AP3D) evaluation, setting Intersection over Union (IoU) threshold to 0.5 and 0.7.

## 4.2. Experimental results

**3D object detection.** We compare the performance with 4 monocular 3D object detection methods: Mono3D [2], OFT-Net [37], MonoGRNet [33] and MF3D [46], which reported their results on the same validation set for the car class. We borrow the presented average precision numbers from their published results. The results are reported in Table 1. The experiments show that our approach outperforms state-of-the-art methods on the easy subset by a small margin while remaining the second best on the moderate subset. This observation aligns with our intuition that visible salient features such as keypoints are crucial to the success of 3D pose estimation. For the moderate and hard images, 2D keypoints are challenging to robustly detect due to the high occlusion level or the low resolution of the instance. We also measure the effect of the reprojection consistency loss on our network performance, observing a positive effect of our loss function.

### 3D bounding box and global orientation estimation.

We follow the experiment presented in [33], evaluating the

quality of the 3D bounding boxes sizes estimation, as well as the orientation in a camera coordinate system. The mean errors of our approach, along with [33, 2], borrowed from their work, are presented in Table 2.

Method	Size (m)			Orientation (rad)
	Height	Width	Length	
Mono3D	0.172	0.103	0.504	0.558
MonoGRNet	<b>0.084</b>	<b>0.084</b>	0.412	0.251
Ours	0.115	0.107	0.516	<u>0.215</u>
Ours (+loss)	<u>0.101</u>	<u>0.091</u>	<b>0.403</b>	<b>0.191</b>

Table 2: **3D bounding box and orientation mean errors:** The best score is in bold, the second best underlined.

Though, the sizes of the 3D bounding boxes do not differ severely among the approaches, due to the estimating the offset from the median bounding box, the orientation estimation results differ significantly. Since we retrieve global orientation via geometric reasoning, learning local orientation from 2D image features, the network provides more accurate predictions, in contrast to obtaining orientation from the regressed 3D bounding box corners.

**Qualitative results.** We provide a qualitative illustration of the network performance in Figure 5, displaying six road scenes with distinct levels of difficulty. In typical cases, our approach produces accurate 3D bounding boxes for all instances, along with the global orientation and 3D location. Remarkably, the truncated objects can also be successfully detected, given that only one pair of keypoints hits the image. Some hard cases, i.e. (e) and (f), primarily consist of objects that are distant, highly occluded or even invisible on the image. We believe such failure cases to be a common limitation of monocular image processing methods.

## 5. Conclusions

In this work, we presented a novel deep learning-based framework for monocular 3D object detection combining well-known detectors with geometric reasoning on keypoints. We proposed to estimate correspondences between the detected 2D keypoints and their 3D counterparts annotated on the surface of 3D CAD models to solve the object localization problem. Results of the experimental evaluation of our approach on the subsets of the KITTI 3D object detection benchmark demonstrate that it outperforms the competing state-of-the-art approaches when the target objects are clearly visible, leading us to hypothesize that dense depth estimation is redundant for 3D detection in some instances. We have demonstrated our multi-task reprojection

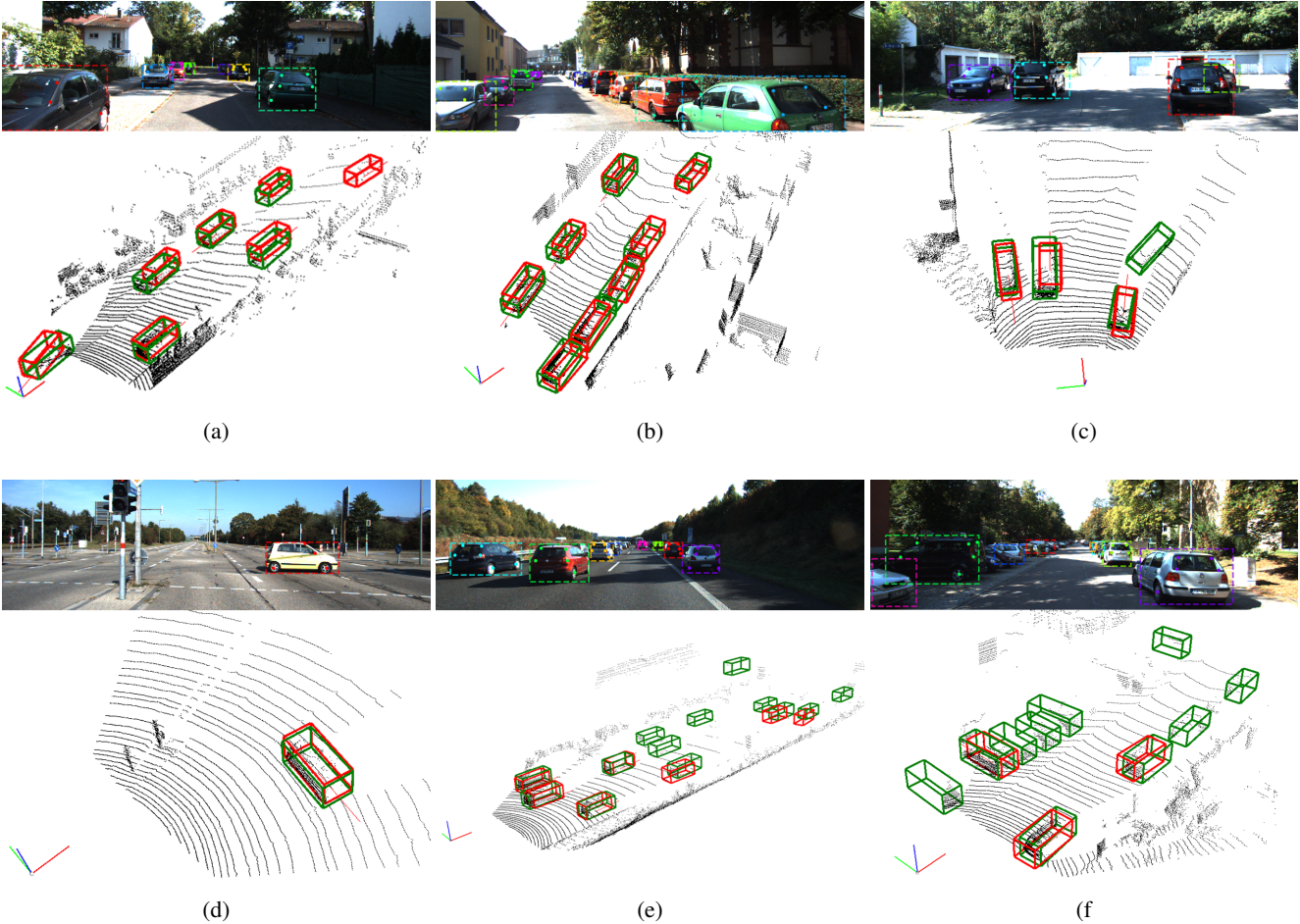


Figure 5: Qualitative results. The upper part of each sub-figure contains 2D detection inference, including 2D bounding and 2D locations of the visible keypoints. Each instance and its keypoints are displayed their distinctive color. The lower part visualizes the 3D point cloud, showing the camera location as the colored XYZ axes. Green and red colors stand for the ground truth and predicted 3D bounding boxes respectively. The scenes were selected to express diversity in complexity and cars positioning w.r.t. the camera.

consistency loss to significantly improve performance, in particular, the orientation of detections.

## Acknowledgement

E. Burnaev and A. Artemov were supported by the Russian Science Foundation under Grant 19-41-04109.

## References

- [1] 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, 2017. [2](#), [3](#)
- [2] X. Chen, K. Kundu, Z. Zhang, H. Ma, S. Fidler, and R. Urtasun. Monocular 3d object detection for autonomous driving. In *CVPR*, 2016. [2](#), [6](#), [7](#)
- [3] X. Chen, H. Ma, J. Wan, B. Li, and T. Xia. Multi-view 3d object detection network for autonomous driving. In *IEEE CVPR*, 2017. [1](#), [2](#)
- [4] J. Dai, Y. Li, K. He, and J. Sun. R-fcn: Object detection via region-based fully convolutional networks. In *Advances in neural information processing systems*, pages 379–387, 2016. [2](#)
- [5] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009. [2](#), [7](#)
- [6] M. Everingham, L. V. Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge, 2010. [2](#), [6](#)
- [7] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013. [2](#)



- [8] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 3354–3361. IEEE, 2012. [2](#)
- [9] C. Godard, O. Mac Aodha, and G. J. Brostow. Unsupervised monocular depth estimation with left-right consistency. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 270–279, 2017. [2](#)
- [10] P. Guerrero, Y. Kleiman, M. Ovsjanikov, and N. J. Mitra. PCPNet: Learning local shape properties from raw point clouds. *Computer Graphics Forum*, 37(2):75–85, 2018. [2](#)
- [11] K. He, G. Gkioxari, P. Dollr, and R. Girshick. Mask r-cnn, 2017. [1](#), [2](#), [3](#), [4](#), [6](#)
- [12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition, 2015. [6](#)
- [13] M. Hejrati and D. Ramanan. Analyzing 3d objects in cluttered images. In *Advances in Neural Information Processing Systems*, pages 593–601, 2012. [1](#)
- [14] B.-S. Hua, M.-K. Tran, and S.-K. Yeung. Pointwise convolutional neural networks, 2017. cite arxiv:1712.05245Comment: 10 pages, 6 figures, 10 tables. Paper accepted to CVPR 2018. [2](#)
- [15] J. Ku, M. Mozifian, J. Lee, A. Harakeh, and S. Waslander. Joint 3d proposal generation and object detection from view aggregation, 2017. [1](#), [2](#)
- [16] A. H. Lang, S. Vora, H. Caesar, L. Zhou, J. Yang, and O. Beijbom. Pointpillars: Fast encoders for object detection from point clouds, 2018. [1](#), [2](#)
- [17] Y. Li, R. Bu, M. Sun, and B. Chen. Pointcnn. *CoRR*, abs/1801.07791, 2018. [2](#)
- [18] J. J. Lim, H. Pirsivash, and A. Torralba. Parsing ikea objects: Fine pose estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2992–2999, 2013. [3](#)
- [19] T.-Y. Lin, P. Dollr, R. Girshick, K. He, B. Hariharan, and S. Belongie. Feature pyramid networks for object detection, 2016. [1](#), [4](#), [6](#)
- [20] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollr. Focal loss for dense object detection, 2017. [1](#), [2](#)
- [21] T.-Y. Lin, M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays, P. Perona, D. Ramanan, C. L. Zitnick, and P. Dollr. Microsoft coco: Common objects in context, 2014. [2](#)
- [22] Y.-L. Lin, V. I. Morariu, W. Hsu, and L. S. Davis. Jointly optimizing 3d model fitting and fine-grained classification. In *European Conference on Computer Vision*, pages 466–480. Springer, 2014. [3](#)
- [23] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia. Path aggregation network for instance segmentation, 2018. [1](#), [2](#)
- [24] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. Ssd: Single shot multibox detector, 2015. [1](#), [2](#)
- [25] P. Marion, P. R. Florence, L. Manuelli, and R. Tedrake. Label fusion: A pipeline for generating ground truth labels for real rgbd data of cluttered scenes. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1–8. IEEE, 2018. [3](#)
- [26] J. Martinez, R. Hossain, J. Romero, and J. J. Little. A simple yet effective baseline for 3d human pose estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2640–2649, 2017. [1](#), [3](#)
- [27] D. Mehta, H. Rhodin, D. Casas, P. Fua, O. Sotnychenko, W. Xu, and C. Theobalt. Monocular 3d human pose estimation in the wild using improved cnn supervision. In *2017 International Conference on 3D Vision (3DV)*, pages 506–516. IEEE, 2017. [1](#), [3](#)
- [28] D. Mehta, S. Sridhar, O. Sotnychenko, H. Rhodin, M. Shafiei, H.-P. Seidel, W. Xu, D. Casas, and C. Theobalt. Vnect: Real-time 3d human pose estimation with a single rgb camera. volume 36, July 2017. [1](#), [3](#)
- [29] A. Mousavian, D. Anguelov, J. Flynn, and J. Kosecka. 3d bounding box estimation using deep learning and geometry. In *CVPR*, 2017. [2](#), [3](#), [5](#)
- [30] C. R. Qi, W. Liu, C. Wu, H. Su, and L. J. Guibas. Frustum pointnets for 3d object detection from rgb-d data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 918–927, 2018. [2](#), [3](#)
- [31] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 652–660, 2017. [2](#)
- [32] C. R. Qi, L. Yi, H. Su, and L. J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in Neural Information Processing Systems*, pages 5099–5108, 2017. [2](#)
- [33] Z. Qin, J. Wang, and Y. Lu. Monogrnet: A geometric reasoning network for monocular 3d object localization. *arXiv preprint arXiv:1811.10247*, 2018. [1](#), [2](#), [7](#)
- [34] V. Ramakrishna, T. Kanade, and Y. Sheikh. Reconstructing 3d human pose from 2d image landmarks. In *European Conference on Computer Vision*, pages 573–586. Springer, 2012. [1](#), [3](#)
- [35] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection, 2015. [1](#), [2](#)
- [36] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks, 2015. [1](#), [2](#), [4](#)
- [37] T. Roddick, A. Kendall, and R. Cipolla. Orthographic feature transform for monocular 3d object detection. *arXiv preprint arXiv:1811.08188*, 2018. [2](#), [7](#)
- [38] N. Snavely, S. M. Seitz, and R. Szeliski. Photo tourism: exploring photo collections in 3d. In *ACM transactions on graphics (TOG)*, volume 25, pages 835–846. ACM, 2006. [3](#)

- [39] S. Suwajanakorn, N. Snavely, J. J. Tompson, and M. Norouzi. Discovery of latent 3d keypoints via end-to-end geometric reasoning. In *Advances in Neural Information Processing Systems*, pages 2063–2074, 2018. 3
- [40] J. J. Tompson, A. Jain, Y. LeCun, and C. Bregler. Joint training of a convolutional network and a graphical model for human pose estimation. In *Advances in neural information processing systems*, pages 1799–1807, 2014. 4
- [41] J. R. Uijlings, K. E. Van De Sande, T. Gevers, and A. W. Smeulders. Selective search for object recognition. *International journal of computer vision*, 104(2):154–171, 2013. 2
- [42] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon. Dynamic graph cnn for learning on point clouds. *CoRR*, abs/1801.07829, 2018. 2
- [43] 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*, pages 365–382. Springer, 2016. 1, 3
- [44] Y. Xiang, W. Choi, Y. Lin, and S. Savarese. Data-driven 3d voxel patterns for object category recognition. In *Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition*, 2015. 6
- [45] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1492–1500, 2017. 2
- [46] B. Xu and Z. Chen. Multi-level fusion based 3d object detection from monocular images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2345–2353, 2018. 1, 2, 7
- [47] D. Xu, D. Anguelov, and A. Jain. Pointfusion: Deep sensor fusion for 3d bounding box estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 244–253, 2018. 2
- [48] Y. Yan, Y. Mao, and B. Li. Second: Sparsely embedded convolutional detection. *Sensors*, 18(10):3337, 2018. 1, 2
- [49] B. Yang, W. Luo, and R. Urtasun. Pixor: Real-time 3d object detection from point clouds. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 1, 2
- [50] M. Zeeshan Zia, M. Stark, and K. Schindler. Explicit occlusion modeling for 3d object class representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3326–3333, 2013. 3
- [51] Y. Zhou and O. Tuzel. Voxynet: End-to-end learning for point cloud based 3d object detection. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 1, 2
- [52] M. Z. Zia, M. Stark, B. Schiele, and K. Schindler. Revisiting 3d geometric models for accurate object shape and pose. In *2011 IEEE International Conference on Computer Vision Workshops (ICCV Workshops)*, pages 569–576. IEEE, 2011. 3
- [53] M. Z. Zia, M. Stark, B. Schiele, and K. Schindler. Detailed 3d representations for object modeling and recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(11):2608–2623, 2013. 3
- [54] C. L. Zitnick and P. Dollár. Edge boxes: Locating object proposals from edges. In *European conference on computer vision*, pages 391–405. Springer, 2014. 2