

# OMNI3D: A Large Benchmark and Model for 3D Object Detection in the Wild

Garrick Brazil Julian Straub Nikhila Ravi  
Justin Johnson Georgia Gkioxari

Meta AI

<https://github.com/facebookresearch/omni3d>

## Abstract

Recognizing scenes and objects in 3D from a single image is a longstanding goal of computer vision with applications in robotics and AR/VR. For 2D recognition, large datasets and scalable solutions have led to unprecedented advances. In 3D, existing benchmarks are small in size and approaches specialize in few object categories and specific domains, *e.g.* urban driving scenes. Motivated by the success of 2D recognition, we revisit the task of 3D object detection by introducing a large benchmark, called OMNI3D. OMNI3D re-purposes and combines existing datasets resulting in 234k images annotated with more than 3 million instances and 97 categories. 3D detection at such scale is challenging due to variations in camera intrinsics and the rich diversity of scene and object types. We propose a model, called Cube R-CNN, designed to generalize across camera and scene types with a unified approach. We show that Cube R-CNN outperforms prior works on the larger OMNI3D and existing benchmarks. Finally, we prove that OMNI3D is a powerful dataset for 3D object recognition, show that it improves single-dataset performance and can accelerate learning on new smaller datasets via pre-training.

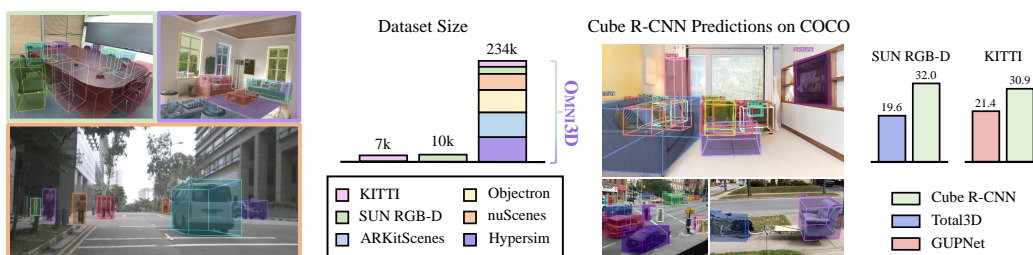


Figure 1: **Left:** We introduce OMNI3D, a benchmark for visual 3D object detection which is larger and more diverse than popular 3D benchmarks. **Right:** We propose Cube R-CNN, which generalizes to unseen images (*e.g.* from COCO [42]) and outperforms prior works on popular 3D benchmarks.

## 1 Introduction

Understanding objects and their properties from single images is a longstanding problem in computer vision with applications in robotics, assistive technology and AR/VR. In the last decade, 2D object recognition [24, 37, 60, 61, 67] has made tremendous advances toward predicting and localizing objects on the 2D image grid with the help of large-scale datasets [22, 42]. However, the world and its objects are three dimensional laid out in 3D space. Perceiving objects in 3D from 2D visual inputs poses new challenges framed by the task of 3D object detection. Here, the goal is to estimate a 3D location and 3D extent of each object in an image in the form of a tight oriented 3D bounding box.

Today 3D object detection is studied under two different lenses: for urban domains in the context of self-driving vehicles [6, 12, 47, 49, 54] or indoor scenes [29, 34, 55, 68]. Despite the problem formulation being shared, methods for urban and indoor scenes share little insights between domains. Often approaches are tailored to work only for the domain in question. For instance, urban methods make assumptions about objects resting on a ground plane and model *only* yaw angles for 3D rotation. Indoor techniques may use a confined depth range (*e.g.* up to 6m in [55]). These assumptions are generally not true for objects and scenes in the real-world. Moreover, the most popular benchmarks for image-based 3D object detection are very small. Indoor SUN RGB-D [66] contains 10k images and the urban KITTI [19] has 7k images; in contrast 2D benchmarks like COCO [42] are  $20\times$  larger.

We address the absence of a general large-scale dataset for 3D object detection by introducing a large and diverse 3D benchmark called OMNI3D. OMNI3D is curated from publicly released datasets, SUN RGB-D [66], ARKitScenes [5], Hypersim [62], Objectron [2], KITTI [19] and nuScenes [8], and comprises 234k images with 3 million objects annotated with 3D boxes across 97 categories including *chair, sofa, laptop, table, cup, shoes, pillow, books, car, person, etc.* As shown in Fig. 1(a), OMNI3D is  $20\times$  larger than existing popular benchmarks used for 3D detection, SUN RGB-D and KITTI. For efficient evaluation on our large dataset, we introduce a new, fast, batched and exact algorithm for intersection-over-union for 3D boxes which is  $450\times$  faster than previous solutions [2]. We empirically prove the impact of OMNI3D as a large-scale dataset and show that it improves single-dataset AP performance by up to 5.3% on urban and 3.8% on indoor benchmarks.

On this new dataset, we design a general and simple 3D object detector, called Cube R-CNN, inspired by key research advances in 2D and 3D recognition of recent years [20, 49, 61, 64], which achieves state-of-the-art results across domains. Cube R-CNN detects all objects and their 3D properties, rotation, and depth end-to-end from a single image of any domain and for many object categories. Attributed to OMNI3D’s diversity, our model shows strong generalization and outperforms prior works for indoor and urban domains with one unified model, as shown in Fig. 1(b). Learning from such diverse data comes with challenges as OMNI3D contains images of highly varying focal lengths which exaggerate scale-depth ambiguity, exemplified in Fig. 4. We remedy this by operating on *virtual depth* which transforms object depth with the same virtual camera intrinsics across the dataset. An added benefit of virtual depth is that it allows the use of data augmentations (*e.g.* image rescaling) during training, which is a critical feature for 2D detection [11, 73], and as we show, also for 3D. Our approach with one unified design outperforms prior state-of-the-art approaches, Total3D [55] by 12.4% IoU<sub>3D</sub> on indoor SUN RGB-D and GUPNet [49] by 9.5% AP<sub>3D</sub> on urban KITTI.

## 2 Related Work

In this work we propose a new benchmark for *general* 3D object detection. We design a unified approach to detect 3D objects from single RGB images by extending 2D object detection with the ability to detect 3D cuboids in an intrinsics-invariant manner.

**2D Object Detection.** Here, methods include two-stage approaches [24, 61] which predict object regions with a region proposal network (RPN) and then refine them via an MLP. Single-stage detectors [41, 44, 60, 67, 78] omit the RPN and predict regions directly from the backbone.

**3D Object Detection.** Monocular 3D object detectors predict 3D cuboids from single input images. There is extensive work in the urban self-driving domain where the *car* class is at the epicenter [12, 17, 21, 23, 28, 43, 46, 51, 54, 59, 64, 65, 71, 77]. CenterNet [78] predicts 3D depth and size from fully-convolutional center features, and is extended by [14, 39, 45, 47–49, 52, 76, 77, 80]. M3D-RPN [6] trains an RPN with 3D anchors, enhanced further by [7, 16, 35, 70, 81]. Others use pseudo depth [3, 13, 50, 56, 72, 74] and explore depth and point-based LiDAR techniques [33, 57]. Similar to ours, [10, 64, 65] add a 3D head, specialized for urban scenes and objects, on two-stage Faster R-CNN. [23, 65] augment their training by synthetically generating depth and box-fitted views, coined as virtual views or depth. In our work, virtual depth aims at addressing *varying focal lengths*.

For indoor scenes, a vast line of work tackles room layout estimation [15, 26, 38, 53]. Huang *et al* [29] predict 3D oriented bounding boxes for indoor objects. Factored3D [68] and 3D-RelNet [34] jointly predict object voxel shapes. Total3D [55] predicts 3D boxes and meshes by additionally training on datasets with annotated 3D shapes. We explore 3D object detection in its general form by tackling *both* domains jointly and with a vocabulary of  $5\times$  more object categories than prior works.

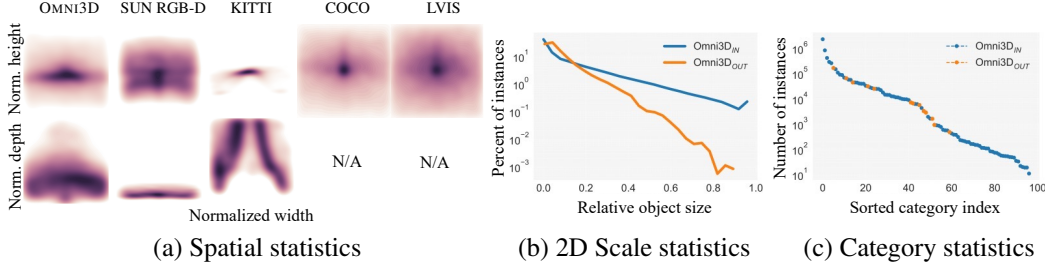


Figure 2: OMNI3D (a) distribution of object centers in normalized image coordinates XY-plane (top) and normalized depth XZ-plane (bottom), (b) relative 2D object scales, and (c) category frequency.

**3D Datasets.** KITTI [19] and SUN RGB-D [66] are popular datasets for 3D object detection on urban and indoor scenes respectively. Since 2019, more datasets have emerged, both for indoor [2, 5, 62] and urban [1, 8, 9, 18, 27, 30, 31]. In isolation, these datasets have unique properties and biases, *e.g.* object and scene types, focal length, coordinate systems, *etc.* In this work, we unify existing representative datasets [2, 5, 8, 19, 62, 66] to form a large-scale benchmark, called OMNI3D, which is  $20\times$  larger than widely-used benchmarks and significantly more diverse. As such, new challenges arise stemming from the increased variance in object rotation, size, layouts, and camera intrinsics.

### 3 The OMNI3D Benchmark

The primary benchmarks for 3D object detection are small (7k images in KITTI [19], 10k in SUN RGB-D [66]), primarily focused on a few categories (*e.g.* 3 on KITTI), and dominated by a single domain. The small size and lack of variance in 3D datasets is a stark difference to 2D counterparts, such as COCO [42] and LVIS [22], which have pioneered progress in 2D recognition. We aim to bridge the gap to 2D by introducing OMNI3D, a large-scale benchmark for image-based 3D object detection consisting of 234k images, 3 million labeled 3D bounding boxes, and 97 object categories. We re-purpose recently released 3D datasets of urban [8, 19], indoor [5, 62, 66], and general [2] scenes to build a diverse large-scale dataset of varying camera intrinsics suitable for general image-based 3D object detection. We analyze OMNI3D and show its rich spatial and semantic properties proving it is visually diverse, similar to 2D data, and highly challenging for 3D as depicted in Fig. 2. We show the value of OMNI3D for the task of 3D object detection with extensive quantitative analysis in Sec. 5.

#### 3.1 Dataset Analysis

We curate OMNI3D from a collection of publicly available data sources, namely SUN RGB-D [66], ARKitScenes [5], Hypersim [62], Objectron [2], KITTI [19] and nuScenes [8]. Each of these datasets are either designed specifically for 3D object detection or contain 3D cuboid annotations which we re-purpose accordingly. We semantically align object categories from each data source via manual association which results in a set of 97 categories. Critically, all cuboids are converted to a *unified 3D camera coordinate system*. There are a total of 234,152 images with more than 3 million annotations. We split the dataset into 175k images for train, 19k for val and 39k for test, consistent with splits from original sources when available, and otherwise free of video sequences overlapping in splits. OMNI3D has images from indoor and outdoor domains. We denote respective subsets as OMNI3D<sub>IN</sub> (SUN RGB-D, Hypersim, ARKit), and OMNI3D<sub>OUT</sub> (KITTI, nuScenes). Note that Objectron, which primarily comprises close-up objects, is used only in the *full* OMNI3D setting.

**Layout statistics.** Fig. 2(a) shows in the top row the distribution of object centers onto the image plane by projecting centroids on the XY-plane, and in the bottom row the distribution of object depths by projecting centroids onto the XZ-plane. We compare OMNI3D to the most popular benchmarks for 3D object detection (SUN RGB-D and KITTI) and 2D object detection (COCO [42] and LVIS [22]). We find that OMNI3D’s spatial distribution has a center bias, similar to 2D datasets. Fig. 2(b) depicts the relative object size distribution, defined as the square root of object area divided by image area. Objects are more likely to be small in size similar to LVIS (Fig. 6c in [22]) suggesting that OMNI3D is also challenging for 2D detection, while objects in OMNI3D<sub>OUT</sub> are noticeably smaller.

The bottom row of Fig. 2(a) normalizes object depth in a  $[0, 20\text{m}]$  range, chosen for visualization and satisfies 88% of object instances; OMNI3D depth ranges as far as 300m. We observe that the

OMNI3D depth distribution is far more diverse than SUN RGB-D and KITTI, which are biased toward near or road-side objects respectively. See Appendix for each data source plot distribution. Fig. 2(a) demonstrates OMNI3D’s rich diversity in spatial distribution and depth which suffers significantly less bias than existing 3D benchmarks and is comparable in complexity to common 2D datasets.

**2D and 3D correlation.** To address the common assumption that urban objects appear on a mostly planar ground, we compute the correlation between normalized 2D  $y$  and 3D  $z$  depth. We find that such attributes are indeed fairly correlated in OMNI3D<sub>OUT</sub> at 0.524, but *significantly* less in OMNI3D<sub>IN</sub> at 0.006 correlation. Similarly, the correlation between relative 2D object size ( $\sqrt{h \cdot w}$ ) and  $z$  is 0.543 and 0.102 respectively. This confirms our intuition that common assumptions in urban scenes do not hold well in other domains, thus making the problem more challenging overall.

**Category statistics.** Fig. 2(c) plots the distribution of instances across the 97 categories of OMNI3D. The long-tail suggests that low-shot object detection in both 2D and 3D will be critical for performance. In this work, we want to focus on large-scale and diverse 3D objects, which is comparably unexplored. We therefore filter and focus on categories that have at least 1000 annotations in 3D. This leaves 50 diverse categories including *chair, sofa, laptop, table, cup, shoes, books, window, car, truck, pedestrian* and many more, with 19 common categories having more than 10k unique instances.

## 4 Method

Our goal is to design a simple and effective model for *general* 3D object detection. Hence, our approach is free of domain or object specific priors. We design our 3D object detection framework as an extension of Faster R-CNN with a 3D object head to predict a cuboid per each detected 2D object. We refer to our method as *Cube R-CNN*. Figure 3 shows an overview of our approach.

### 4.1 Cube R-CNN

Our model builds on Faster R-CNN [61], an end-to-end region-based object detection framework. Faster-RCNN consists of a backbone network, commonly a CNN, which embeds the input image into a higher-dimensional feature space. A region proposal network (RPN) predicts regions of interest (RoIs) representing object candidates in the image. A 2D box head inputs the backbone feature map and processes each RoI to predict a category and a more accurate 2D bounding box. Faster R-CNN can be easily extended to tackle more tasks by adding task-specific heads *e.g.* Mask R-CNN [24] adds a mask head to additionally output object silhouettes.

For the task of 3D object detection, we extend Faster R-CNN by introducing a 3D detection head which predicts a 3D cuboid for each detected 2D object. Cube R-CNN extends Faster R-CNN in three ways: (a) we replace the binary classifier in RPN which predicts region *objectness* with a regressor that predicts *IoUness*, (b) we introduce a *cube head* which estimates the parameters to define a 3D cuboid for each detected object (similar in concept to [64]), and (c) we define a new training objective which incorporates a *virtual depth* for the task of 3D object detection.

**IoUness.** The role of RPN is two-fold: (a) it proposes RoIs by regressing 2D box coordinates from pre-computed anchors and (b) it classifies regions as *object or not* (objectness). This is sensible in exhaustively labeled datasets where it can be reliably assessed if a given region contains an object. However, OMNI3D combines many data sources with no guarantee that all instances of all classes are labeled. We overcome this by replacing objectness with *IoUness*, applied only to foreground. Similar to [32], we replace objectness in the RPN with a regressor to predict IoU between an RoI and a ground truth. Let  $o$  be the predicted IoU for a RoI and  $\hat{o}$  be the 2D IoU between the region and its ground truth; we apply a binary cross-entropy (CE) loss  $\mathcal{L}_{\text{IoUness}} = \ell_{\text{CE}}(o, \hat{o})$ . We train on regions whose IoU exceeds 0.05 with a ground truth in order to learn *IoUness* from a wide range of region overlaps. Thus, the RPN training objective becomes  $\mathcal{L}_{\text{RPN}} = \hat{o} \cdot (\mathcal{L}_{\text{IoUness}} + \mathcal{L}_{\text{reg}})$ , where  $\mathcal{L}_{\text{reg}}$  is the 2D box regression loss from [61]. The loss is weighted by  $\hat{o}$  to prefer candidates close to true objects.

**Cube Head.** We extend Faster R-CNN with a new head, called cube head, to predict a 3D cuboid for each detected 2D object. The cube head inputs  $7 \times 7$  feature maps pooled from the backbone for each predicted region and feeds them to 2 fully-connected (FC) layers with 1024 hidden dimensions. Similar to Faster R-CNN, all 3D estimations in the cube head are category-specific.



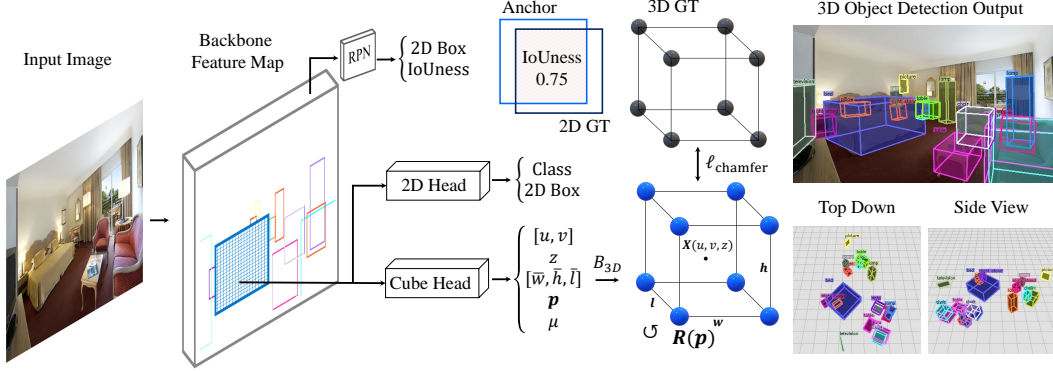


Figure 3: Cube R-CNN takes as input an RGB image, detects all objects in 2D and predicts their 3D cuboids  $B_{3D}$ . During training, the 3D cuboid corners are compared against a 3D ground truth.

The cube head represents a 3D cuboid with 13 parameters each predicted by a final FC layer:

- $[u, v]$  represent the projected 3D center on the image plane relative to the 2D RoI
- $z \in \mathbb{R}_+$  is the object's center depth in meters transformed from virtual depth  $z_v$  (see below)
- $[\bar{w}, \bar{h}, \bar{l}] \in \mathbb{R}_+^3$  are the log-normalized physical box dimensions in meters
- $\mathbf{p} \in \mathbb{R}^6$  represents a continuous 6D [79] allocentric rotation
- $\mu \in \mathbb{R}$  is the predicted 3D uncertainty

The above parameters form the final 3D box in camera view coordinates for each detected object. The object's 3D center  $\mathbf{X}$  is estimated from the predicted 2D projected center  $[u, v]$  and depth  $z$  via

$$\mathbf{X}(u, v, z) = \left( \frac{z}{f_x}(r_x + ur_w - p_x), \frac{z}{f_y}(r_y + vr_h - p_y), z \right) \quad (1)$$

where  $[r_x, r_y, r_w, r_h]$  is the object's 2D box,  $(f_x, f_y)$  are the camera's known focal lengths and  $(p_x, p_y)$  the principal point. The 3D box dimensions  $\mathbf{d}$  are derived from  $[\bar{w}, \bar{h}, \bar{l}]$  which are log-normalized with category-specific pre-computed means  $(w_0, h_0, l_0)$  for width, height and length respectively, which are combined then arranged into a diagonal scaling matrix via

$$\mathbf{d}(\bar{w}, \bar{h}, \bar{l}) = \text{diag}(\exp(\bar{w})w_0, \exp(\bar{h})h_0, \exp(\bar{l})l_0) \quad (2)$$

Finally, we derive the rotation  $\mathbf{R}(\mathbf{p})$  of the object as a  $3 \times 3$  rotation matrix based on a 6D parameterization (2 directional vectors) of  $\mathbf{p}$  following [79] which is converted from allocentric to egocentric rotation similar to [36], defined formally in Appendix. The final 3D cuboid, defined by 8 corners, is

$$B_{3D}(u, v, z, \bar{w}, \bar{h}, \bar{l}, \mathbf{p}) = \mathbf{R}(\mathbf{p}) \mathbf{d}(\bar{w}, \bar{h}, \bar{l}) B_{\text{unit}} + \mathbf{X}(u, v, z) \quad (3)$$

where  $B_{\text{unit}}$  are the 8 corners of an axis-aligned unit cube centered at  $(0, 0, 0)$ . Lastly,  $\mu$  denotes a learned 3D uncertainty, which is mapped to a confidence at inference then joined with the classification score  $s$  from the 2D box head to form the final score for the prediction, as  $\sqrt{s \cdot \exp(-\mu)}$ .

**Training objective.** Our training objective consists of 2D losses from the RPN and 2D box head and 3D losses from the cube head. The 2D losses follow [61] with a deviation for the RPN loss which replaces objectness with IoUness, as described above. The cube head introduces the 3D losses. We begin by comparing each predicted 3D cuboid with its matched ground truth via a chamfer loss, treating the 8 box corners of the 3D boxes as point clouds, namely  $\mathcal{L}_{3D}^{\text{all}} = \ell_{\text{chamfer}}(B_{3D}, B_{3D}^{\text{gt}})$ . Note that  $\mathcal{L}_{3D}^{\text{all}}$  entangles all 3D variables via the box predictor  $B_{3D}$  (Eq. 3), such that that errors in variables may be ambiguous from one another. Thus, following [64], we isolate each variable group with separate disentangled losses. The disentangled loss for each variable group substitutes all but its variables with the ground truth from Eq. 3 to create a pseudo box prediction. For example, the disentangled loss for the projected 3D center  $[u, v]$  would produce a 3D box with all but  $(u, v)$  replaced with the *true* values and then compare against the ground truth box, via

$$\mathcal{L}_{3D}^{(u,v)} = \|B_{3D}(u, v, z^{\text{gt}}, \bar{w}^{\text{gt}}, \bar{h}^{\text{gt}}, \bar{l}^{\text{gt}}, \mathbf{p}^{\text{gt}}) - B_{3D}^{\text{gt}}\|_1 \quad (4)$$

We use an  $L_1$  loss for  $\mathcal{L}_{3D}^{(u,v)}$ ,  $\mathcal{L}_{3D}^{(z)}$  and  $\mathcal{L}_{3D}^{(\bar{w}, \bar{h}, \bar{l})}$  and chamfer for  $\mathcal{L}_{3D}^{\mathbf{p}}$  to account for cuboid symmetry such that rotation matrices of Euler angles modulo  $\pi$  produce the same non-canonical 3D cuboid.

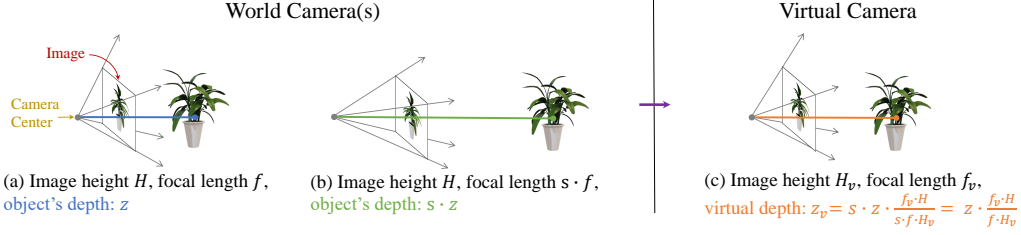


Figure 4: Varying camera intrinsics exaggerate scale-depth ambiguity as the same object at two different depths can project to the same image, as shown in (a) and (b). We address this by introducing a *virtual camera system* which is invariant to intrinsics and transforms the object’s depth to a *virtual depth*  $z_v$  such that the effective image size  $H_v$  and focal length  $f_v$  are consistent across images.

Losses in box coordinates have natural advantages over losses which directly compare variables to ground truths in that their gradients are appropriately weighted by the error as shown in [64].

The collective training objective of Cube R-CNN is given by,

$$\mathcal{L} = \mathcal{L}_{\text{RPN}} + \mathcal{L}_{2\text{D}} + \sqrt{2} \cdot \exp(-\mu) \cdot (\mathcal{L}_{3\text{D}}^{\text{all}} + \mathcal{L}_{3\text{D}}^{(u,v)} + \mathcal{L}_{3\text{D}}^{(z)} + \mathcal{L}_{3\text{D}}^{(\bar{w}, \bar{h}, \bar{l})} + \mathcal{L}_{3\text{D}}^{\text{P}}) + \mu \quad (5)$$

where  $\mathcal{L}_{\text{RPN}}$  is the RPN loss, described above, and  $\mathcal{L}_{2\text{D}}$  is the loss from the 2D box head from [61]. The 3D losses are weighted by the predicted 3D uncertainty (inspired by [49]), such that the model may trade a penalty to reduce the 3D loss when uncertain. In practice,  $\mu$  is helpful for both improving the 3D rating of a cuboids at inference and reducing the loss of hard samples during training.

## 4.2 Virtual Depth

A critical part of 3D object detection is predicting an object’s depth from the camera in metric units. Estimating depth from only visual cues requires an implicit mapping of 2D pixels to 3D distances, which is more ambiguous if camera intrinsics vary. Prior works are able to ignore this as they primarily train on images from a single sensor. OMNI3D contains images from many sensors and thus demonstrates large variations in camera intrinsics. We design Cube R-CNN to be robust to intrinsics by predicting the object’s *virtual depth*  $z_v$ , which projects the metric depth  $z$  to a virtual invariant camera space. As depicted in Fig. 4, virtual depth scales the depth using the (known) camera intrinsics such that the effective image size and focal length are *consistent across the dataset*.

*Definition:* Virtual depth is designed to maintain an *effective* image height  $H_v$  and focal length  $f_v$ , which are both hyperparameters we refer to as virtual height and virtual focal length, respectively. Let  $z$  be the true metric depth of an object from an image with focal length  $f$  and image height  $H$ . We define the object’s virtual depth as  $z_v = z \cdot \frac{f_v}{f} \cdot \frac{H}{H_v}$ . The derivation can be found in the Appendix.

Invariance to camera intrinsics using virtual depth *also* enables scale augmentations during training, since  $H$  can vary without harming the consistency of  $z_v$ . Data augmentations from image resizing tend to be critical for 2D models but are not commonly used in 3D methods [6, 47, 64] since if unaccounted for they may increase scale-depth ambiguity. With virtual depth, such ambiguities are lessened which therefore enables powerful data augmentation strategies in training. We justify the effects of *virtual depth*, dealing with varying intrinsics and enabling scale augmentations, in Sec. 5.

## 5 Experiments

We evaluate single-image 3D object detection on OMNI3D. We show our method’s superiority using a single model design by comparing to prior best models on existing benchmarks and OMNI3D. Finally, we prove the effectiveness of OMNI3D as a large-scale 3D object detection benchmark by showing comprehensive cross-dataset generalization and its impact for pre-training.

**Implementation details.** We implement Cube R-CNN using Detectron2 [73] and PyTorch3D [58]. We use DLA-34 [75] pretrained on ImageNet [63] with a FPN [40] as our backbone. We train for 128 epochs with a batch size of 192 images distributed across 48 V100s. We use SGD with a learning rate of 0.12 which decays by a factor of 10 after 60% and 80% of training. During training, we use random data augmentation of horizontal flipping and scales  $\in [0.50, 1.25]$ , enabled by virtual depth. Virtual parameters are  $f_v = H_v = 512$ . All source code and models will be made publicly available.

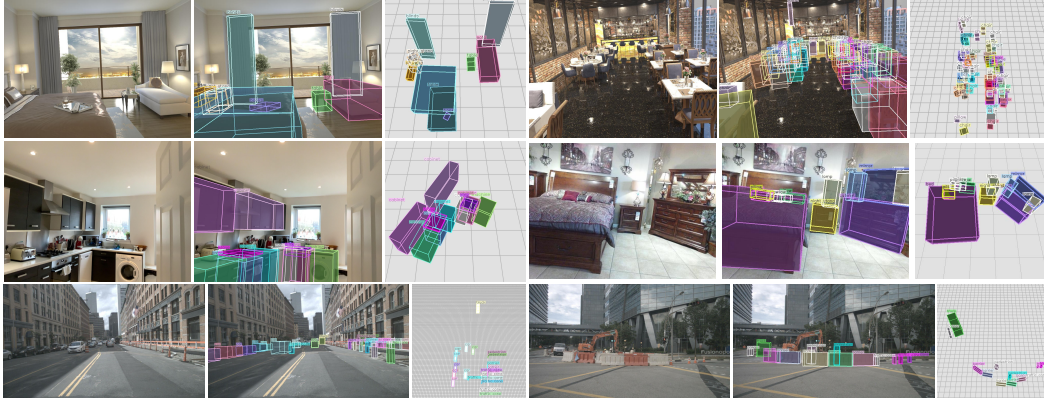


Figure 5: Cube R-CNN predictions on OMNI3D test. For each example, we show the input image, the 3D predictions overlaid on the image and a top view. More examples are given in the Appendix.

**Metric.** Following 2D object detection, we define average-precision (AP) as our core 3D object detection metric. Predictions are matched to ground truth by measuring their overlap using  $\text{IoU}_{3D}$  which computes the intersection-over-union of 3D cuboids, replacing  $\text{IoU}_{2D}$  in 2D. We contribute a novel, fast, and exact algorithm to compute  $\text{IoU}_{3D}$ , described below. We follow COCO [42] and compute a *mean*  $\text{AP}_{3D}$  across all 50 categories in OMNI3D and over a range of  $\text{IoU}_{3D}$  thresholds  $\tau \in [0.05, 0.10, \dots, 0.50]$ . The range of  $\tau$  is more relaxed than in 2D to account for the new dimension. We note that general 3D objects may be occluded by other objects or truncated by the image border, and can be arbitrarily small. Following [19], during evaluation we treat objects with high occlusion or truncation ( $> 66\%$ ), and tiny objects ( $< 6.25\%$  image height) as ignored. We also denote  $\text{AP}_{3D}$  at varying levels of depth  $d$  as near:  $0 < d \leq 10m$ , medium:  $10m < d \leq 35m$ , far:  $35m < d \leq \infty$ .

**Fast  $\text{IoU}_{3D}$ .**  $\text{IoU}_{3D}$  compares two 3D cuboids by computing their intersection-over-union. Images usually have many ground truths and produce several predictions so  $\text{IoU}_{3D}$  computations need to be fast and accurate. Prior implementations [19, 66] approximate  $\text{IoU}_{3D}$  by projecting 3D boxes on a ground plane and multiply the top-view 2D intersection with the box heights to compute a 3D volume. Such approximations become notably inaccurate when objects are not on a planar ground or have arbitrary orientation (*e.g.* nonzero pitch or roll). Objectron [2] provides an unbatched exact solution, but is reliant on external libraries [4, 69] implemented only in C++. We implement a new, fast and exact algorithm which directly finds the intersecting shape by representing cuboids as triangular meshes and finds face intersections. Our algorithm is self-contained, batched with C++ and CUDA support. Our implementation is  $90\times$  faster than Objectron in C++, and  $450\times$  faster in CUDA. As a result, evaluation on the large OMNI3D takes only a few seconds instead of several hours.

## 5.1 Model Performance

We first ablate the design choices of our Cube R-CNN. Then, we compare our approach to prior state-of-the-art methods on existing benchmarks and OMNI3D and show that it performs superior or on par with prior works which are generally specialized for their respective domains. Our method, Cube R-CNN, uses a single unified design to tackle general 3D object detection across domains. Figure 5 shows qualitative predictions on the OMNI3D test set.

**Ablations.** Table 1 ablates the features of our proposed Cube R-CNN on OMNI3D. We report  $\text{AP}_{3D}$  and additional metrics at single IoU thresholds (0.25 and 0.50), as well as performance for near, medium and far objects. We further train and evaluate on domain-specific subsets,  $\text{OMNI3D}_{\text{IN}}$  and  $\text{OMNI3D}_{\text{OUT}}$  ( $\text{AP}_{3D}^{\text{IN}}$  and  $\text{AP}_{3D}^{\text{OUT}}$ , respectively) to show how ablation trends hold for single domain scenarios. From Table 1 we draw a few standout observations:

*Virtual depth* is effective and improves  $\text{AP}_{3D}$  by +6%, most noticable in the full OMNI3D which has the largest variances in camera intrinsics. Enabled by virtual depth, scale augmentations during training increase  $\text{AP}_{3D}$  by +3.1% on OMNI3D, +2.7% on  $\text{OMNI3D}_{\text{IN}}$  and +3.8% on  $\text{OMNI3D}_{\text{OUT}}$ . We find scale augmentation controlled *without* virtual depth to be harmful by  $-2.5\%$  (rows 6 - 7).

Cube R-CNN	AP <sub>3D</sub>	AP <sub>3D</sub> <sup>25</sup>	AP <sub>3D</sub> <sup>50</sup>	AP <sub>3D</sub> <sup>near</sup>	AP <sub>3D</sub> <sup>med</sup>	AP <sub>3D</sub> <sup>far</sup>	AP <sub>3D</sub> <sup>IN</sup>	AP <sub>3D</sub> <sup>OUT</sup>
w/o disentangled	22.7	24.3	9.3	27.9	11.3	8.1	14.7	31.5
w/o IoUness	22.2	23.6	8.8	26.4	11.1	8.3	14.3	31.0
w/o $\mathcal{L}_{3D}^{all}$	20.2	21.8	6.4	26.4	<b>12.1</b>	7.4	13.8	29.1
w/o scale aug.	20.2	21.5	8.0	23.5	9.8	6.8	12.3	28.1
w/o scale l virtual	19.8	21.2	7.5	23.4	8.6	5.7	12.2	26.0
w/o virtual depth	17.3	18.4	4.4	22.7	7.9	7.1	13.2	30.3
w/o uncertainty $\mu$	17.2	18.3	5.4	20.5	10.8	7.0	9.1	25.4
ours	<b>23.3</b>	<b>24.9</b>	<b>9.5</b>	<b>27.9</b>	<b>12.1</b>	<b>8.5</b>	<b>15.0</b>	<b>31.9</b>

Table 1: Cube R-CNN ablations. We report AP<sub>3D</sub> on OMNI3D, at IoU thresholds 0.25 and 0.50, and for near, medium and far objects. We further report AP<sub>3D</sub> on OMNI3D<sub>IN</sub> and OMNI3D<sub>OUT</sub> subsets.

Method	$f$	KITTI			nuScenes		OMNI3D <sub>OUT</sub>		
		AP <sub>3D</sub> <sup>KIT</sup>	AP <sub>3D</sub> <sup>NU</sup>	AP <sub>car</sub> <sup>70</sup>	AP <sub>3D</sub> <sup>KIT</sup>	AP <sub>3D</sub> <sup>NU</sup>	AP <sub>3D</sub> <sup>KIT</sup>	AP <sub>3D</sub> <sup>NU</sup>	AP <sub>3D</sub> <sup>OUT</sup>
M3D-RPN [6]	✗	18.4	0.8	11.4	2.2	18.9	10.4	17.9	13.7
M3D-RPN [6]++	✓	19.0	10.7	9.3	10.6	19.6	16.2	20.4	17.0
GUPNet [49]	✓	21.4	8.5	14.4	17.3	23.1	24.5	20.5	19.9
Cube R-CNN (ours)	✓	<b>30.9</b>	<b>12.7</b>	<b>14.7</b>	<b>20.2</b>	<b>33.0</b>	<b>36.0</b>	<b>32.7</b>	<b>31.9</b>

Table 2: We compare to GUPNet [49], M3D-RPN [6] and its extension with virtual depth and 6D allocentric pose, M3D-RPN++. We report AP<sub>3D</sub>, KITTI AP<sub>car</sub><sup>70</sup> of [19, 64] for val<sub>1</sub> on cars with IoU threshold of 0.7, and zero-shot generalization.  $f$  denotes if a method handles variance in focal length.

*3D uncertainty* boosts performance by about +6% on all OMNI3D subsets. Intuitively, it serves as a measure of confidence for the model’s 3D predictions, and removing it means that the model would rely only on its 2D classification to score cuboids. If uncertainty is used only to scale samples in Eq. 5, but not at inference then AP<sub>3D</sub> still improves but by +1.6%. Table 1 also shows that the entangled  $\mathcal{L}_{3D}^{all}$  loss from Eq. 5 boosts AP<sub>3D</sub> by +3.1% while disentangled losses contribute less (+0.6%). Replacing *objectness* with *IoUness* in the RPN head improves AP<sub>3D</sub> on OMNI3D by +1.1%.

**Comparison with other methods.** We compare Cube R-CNN with state-of-the-art approaches. We choose representative methods designed for the urban and indoor domains, and evaluate on our proposed OMNI3D benchmark and single-dataset benchmarks, KITTI and SUN RGB-D.

*Urban baselines.* Table 2 compares Cube R-CNN to baselines M3D-RPN [6], a single-shot 3D anchor approach, and GUPNet [49], which uses 2D-3D geometry and camera focal length to derive depth. We train M3D-RPN and GUPNet using their publicly available code and report the best performance across 3 runs. For a stronger comparison with [6], we extend it using our virtual depth and the 6D [79] allocentric pose, referred to as M3D-RPN++. We train and evaluate on three datasets: KITTI (columns 3-5), nuScenes (columns 6-7) and OMNI3D<sub>OUT</sub> (columns 8-10) using our defined data splits for trainval and test. For each dataset, we report performance on the test set as well as zero-shot generalization to the other dataset (shown with blue). For OMNI3D<sub>OUT</sub>, we report performance on its test set, AP<sub>3D</sub><sup>OUT</sup>, and its subparts, AP<sub>3D</sub><sup>KIT</sup> and AP<sub>3D</sub><sup>NU</sup>. All methods are trained with the same backbone and training recipe on the 5 KITTI categories, 9 nuScenes categories and 11 OMNI3D<sub>OUT</sub> categories. On KITTI, we additionally report the commonly used, but approximate, KITTI val<sub>1</sub> metric [19, 64] on the *car* category with the strict 0.7 IoU threshold, denoted as AP<sub>car</sub><sup>70</sup>.

From Table 2 we observe that our Cube R-CNN outperforms competitive methods across all datasets on our robust AP<sub>3D</sub> metric by a large margin and performs comparably on the more volatile AP<sub>car</sub><sup>70</sup>. We observe that our method shows the strongest zero-shot generalization performance on both KITTI and nuScenes. M3D-RPN fails to generalize without virtual depth and succeeds with it in M3D-RPN++, which increases its zero-shot generalization to about 10% for either direction of KITTI  $\leftrightarrow$  nuScenes. This shows that virtual depth is effective across methods. Lastly, training Cube R-CNN on the much larger OMNI3D<sub>OUT</sub> greatly impacts its 3D performance, *e.g.* from 30.9% to 36.0% for KITTI.

*Indoor baselines.* Table 3 compares Cube R-CNN to Total3D [55], a state-of-the-art method on SUN RGB-D. We use Total3D’s public model which *requires* 2D object boxes and categories as input. For a fair comparison, we use ground truth 2D detections as input to both Total3D and our Cube R-CNN, each trained using a ResNet34 [25] backbone. We compare Cube R-CNN from two respects, firstly trained only on SUN RGB-D which is identical to Total3D’s setting, and secondly trained on OMNI3D<sub>IN</sub> which subsumes SUN RGB-D. We report the mean IoU<sub>3D</sub> for 12 categories in Table 3. Our model outperforms Total3D by +12.4% when trained on the same training set and increases the performance gap to +13.6% when trained on the large-scale OMNI3D<sub>IN</sub> dataset.



Method	Trained on	bed	books	cabinet	chair	counter	desk	lamp	pillow	fridge	sofa	sink	table	avg.
Total3D [55]	SUN RGB-D	33.6	6.9	13.1	24.2	16.9	23.1	10.5	11.7	19.1	30.1	18.8	27.7	19.6
Cube R-CNN	SUN RGB-D	49.5	<b>10.4</b>	26.6	39.9	34.4	<b>35.7</b>	15.6	<b>20.4</b>	<b>34.7</b>	46.0	<b>31.9</b>	39.2	32.0
Cube R-CNN	OMNI3D <sub>IN</sub>	<b>50.1</b>	10.2	<b>29.0</b>	<b>41.6</b>	<b>37.5</b>	35.6	<b>17.3</b>	<b>20.4</b>	34.6	<b>50.0</b>	31.8	<b>40.7</b>	<b>33.2</b>

Table 3: Comparison to Total3D [55] on SUN RGB-D test. We use oracle 2D detections fairly for all methods and report IoU<sub>3D</sub>. Cube R-CNN outperforms Total3D by +12.4% with equivalent training.

Trained on	AP <sub>3D</sub> <sup>HYP</sup>	AP <sub>3D</sub> <sup>SUN</sup>	AP <sub>3D</sub> <sup>AR</sup>
Hypersim	15.2	9.5	7.5
SUN	5.8	34.7	13.1
ARKit	5.9	14.2	38.6
OMNI3D	<b>19.0</b>	32.6	38.2
OMNI3D <sub>IN</sub>	17.8	<b>35.4</b>	<b>41.2</b>

Trained on	AP <sub>3D</sub> <sup>KIT</sup>	AP <sub>3D</sub> <sup>NU</sup>
KITTI	37.1	12.7
nuScenes	20.2	38.6
OMNI3D	37.8	35.8
OMNI3D <sub>OUT</sub>	<b>42.4</b>	<b>39.0</b>

Table 4: Cross-dataset performance.

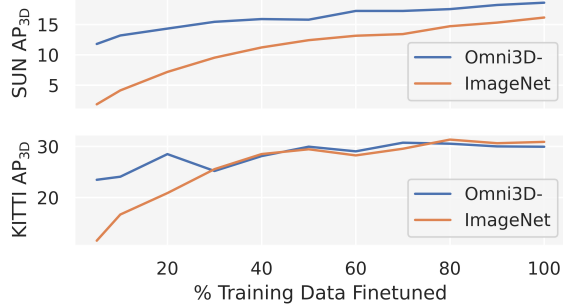


Figure 6: Pre-training on OMNI3D- vs. ImageNet.

## 5.2 The Impact of the OMNI3D Dataset

Sec. 5.1 analyzes the performance of our Cube R-CNN and proves that we outperform prior methods. Now, we turn to OMNI3D and its impact as a large-scale dataset. We show two use cases of OMNI3D: (a) a universal 3D dataset which integrates smaller ones, and (b) a pre-training dataset.

**OMNI3D as a universal dataset.** We treat OMNI3D as a dataset which integrates smaller single ones and show the impact of OMNI3D on each one. To this end, we train Cube R-CNN on OMNI3D and compare to single-dataset training. Table 4 shows the performance for the indoor (top) and urban (bottom) domain reported on the category intersection, 14 for indoor and 3 for outdoor, to ensure AP<sub>3D</sub> comparability. Training on OMNI3D and its domain-specific subsets, OMNI3D<sub>IN</sub> and OMNI3D<sub>OUT</sub>, results in higher performance compared to single-dataset training, signifying that our large-scale OMNI3D generalizes better and should be preferred over single dataset training. Notably, ARKit sees a +2.6% boost and KITTI +5.3%. Except for Hypersim, the domain-specific subsets tend to perform better on their domain, which is not surprising given their distinct properties (Fig. 2).

**OMNI3D as a pre-training dataset.** Next, we demonstrate the utility of OMNI3D for pre-training. In this setting, an unseen dataset is finetuned using a Cube R-CNN model pre-trained from OMNI3D. The motivation is to determine how a large-scale 3D dataset could accelerate low-shot learning with minimum need for costly 3D annotations on a new dataset. We choose SUN RGB-D and KITTI as our *unseen* target datasets given their popularity and small size. We pre-train Cube R-CNN on OMNI3D-, which removes the target datasets from OMNI3D, and then subsequently finetune the models using percentages of their available training data. The performance curves in Fig. 6 show the model quickly gains its final performance even at a small fraction of the training data when pretrained on OMNI3D- compared to ImageNet. When pre-trained on OMNI3D- a model finetuned from only 5% of either target achieves more than 70% of their upper-bound performance.

## 6 Conclusion

We revisit the task of image-based 3D object detection by proposing a large benchmark, OMNI3D, with many categories from diverse domains and with varying camera intrinsics. For OMNI3D, we propose a simple yet effective model, Cube R-CNN, and an exact implementation for IoU for fast evaluation. Future work on this new and exciting benchmark includes few-shot learning on the 47 categories of OMNI3D with few annotations which we do not tackle. Our quantitative analysis shows that 3D object detection is challenging, especially for far away objects. **Ethics:** We explore technical advances for visual 3D recognition, which on the surface are neutral from an ethics perspective. We acknowledge that the used public datasets have a pedestrian category which has ethical consideration, but are otherwise far more dominated by general inanimate objects, chairs, lamps, cars, books *etc.* However, we remain committed to be wary of any harmful consequences stemming from our work.



## Appendix

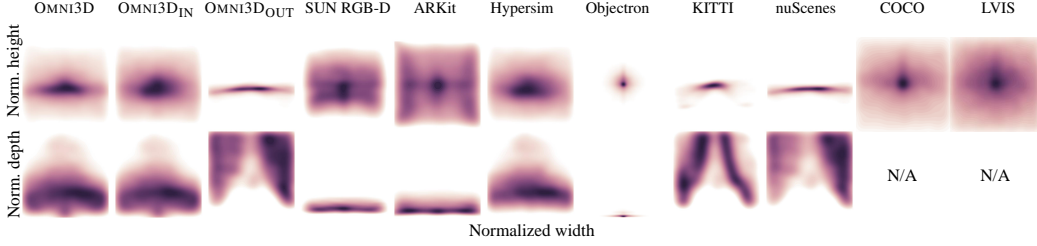


Figure 7: We show the distribution of object centers in normalized image coordinates projected onto the XY-plane (top) and normalized depth projected onto the topview XZ-plane (bottom) across each subset of OMNI3D, individual 3D datasets KITTI [19], SUN RGB-D [66], nuScenes [8], ARKit [5], and Hypersim [62], and large-scale 2D datasets COCO [42] and LVIS [22].

### A1 Dataset Details

In this section, we give more details of OMNI3D, its sources and subsets. We will release the OMNI3D annotations along with our models.

**Sources.** For each individual dataset used in OMNI3D, we use their official train, val, and test set for any datasets which have all annotations released. If no public test set is available then we use their validation set *as* our test set. Whenever necessary we further split the remaining training set in 10:1 ratio by sequences in order form train and val sets. The resultant image make up of OMNI3D is:

- KITTI [19] has 7481 images (3321 train, 391 val, 3769 test).
- SUN RGB-D [66] has 10335 images (4929 train, 356 val, 5050 test).
- nuScenes [8] has 34149 images (26215 train, 1915 val, 6019 test).
- Objectron [2] has 46644 images (33519 train, 3811 val, and 9314 test).
- ARKitScenes [5] has 60924 images (48046 train, 5268 val, and 7610 test).
- Hypersim [62] has 74619 images (59543 train, 7386 val, and 7690 test).

**Coordinate system.** We define our unified 3D coordinate system for all labels with the camera center being the origin and +x facing right, +y facing down, +z inward [19]. Object pose is relative to an initial object with its bottom-face normal aligned to +y and its front-face aligned to +x (*e.g.* upright and facing to the right). All images have known camera calibration matrices with input resolutions varying from 370 to 1920 and diverse focal lengths from 518 to 1708 in pixels. Each object label contains a category label, a 2D bounding box, a 3D centroid in camera space meters, a  $3 \times 3$  matrix defining the object to camera rotation, and the physical dimensions (width, height, length) in meters.

**Spatial Statistics** Following Section 3, we provide the spatial statistics for the individual data sources of KITTI [19], SUN RGB-D [66], nuScenes [8], ARKit [5], and Hypersim [62], as well as OMNI3D<sub>IN</sub> and OMNI3D<sub>OUT</sub> in Figure 7. As in the main paper we observe that the indoor domain data, with the exception of Hypersim, have bias for close objects. Moreover, outdoor data tend to be spatially biased with projected centroids along diagonal ground planes while indoor is more central.

### A2 Model Details

In this section, we provide additional details on our Cube R-CNN model pertaining to its 3D bounding box allocentric rotation (of Section 4.1) and the derivation of virtual depth (of Section 4.2).

#### A2.1 3D Box Rotation

Our 3D bounding box object rotation is predicted in the form of a 6D continuous parameter, which is shown to be well suited for neural networks to regress compared to other forms of rotation in [79].

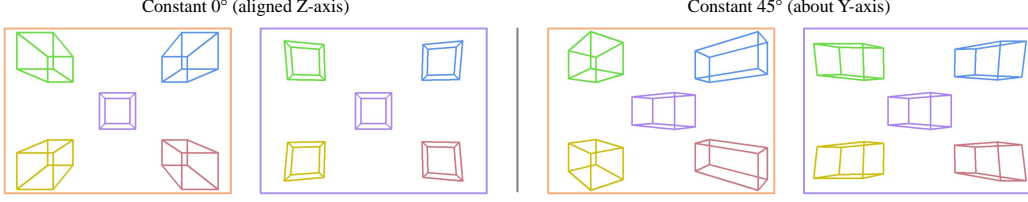


Figure 8: We show **egocentric** and **allocentric** representations under constant rotations. Despite being identical rotations, the egocentric representation appears visually different as its spatial location changes, unlike allocentric which maintains a consistent appearance at all spatial locations.

Let our predicted rotation  $\mathbf{p}$  be split into two directional vectors  $\mathbf{p}_{1-2} \in \mathbb{R}^3$  and  $\mathbf{r}_{1-3} \in \mathbb{R}^3$  be the columns of a  $3 \times 3$  rotational matrix  $\mathbf{R}_a$ . Then  $\mathbf{p}$  is mapped to  $\mathbf{R}_a$  using  $\mathbf{r}_1 = \text{norm}(\mathbf{p}_1)$ ,  $\mathbf{r}_2 = \text{norm}(\mathbf{p}_2 - (\mathbf{r}_1 \cdot \mathbf{p}_2)\mathbf{r}_1)$ , and  $\mathbf{r}_3 = \mathbf{r}_1 \times \mathbf{r}_2$ , where  $\cdot$ ,  $\times$  denote dot and cross product respectively.

$\mathbf{R}_a$  is estimated in *allocentric* form similar to [36]. Let  $f_x, f_y, p_x, p_y$  be the known camera intrinsics,  $u, v$  the predicted 2D projected center as in Section 4.1, and  $a = [0, 0, 1]$  be the camera’s principal axis. Then  $o = \text{norm}([\frac{u-p_x}{f_x}, \frac{v-p_y}{f_y}, 1])$  is a ray pointing from the camera to  $u, v$  with angle  $\alpha = \text{acos}(o)$ . Using standard practices of axis angle representations with an axis denoted as  $o \times a$  and angle  $\alpha$ , we compute a matrix  $\mathbf{M} \in \mathbb{R}^{3 \times 3}$ , which helps form the final *egocentric* rotation matrix  $\mathbf{R} = \mathbf{M} \cdot \mathbf{R}_a$ .

We provide examples of 3D bounding boxes at constant egocentric or allocentric rotations in Figure 8. The allocentric rotation is more aligned to the visual 2D evidence, whereas egocentric rotation entangles relative position into the prediction. In other words, identical egocentric rotations may look very different when viewed from varying spatial locations, which is *not* true for allocentric.

## A2.2 Virtual Depth

In Section 4.2, we propose a virtual depth transformation in order to help our Cube R-CNN model handle varying input image resolutions and camera intrinsics. The motivation of estimating a *virtual* depth instead of *metric* depth is to keep the effective image size and focal length consistent in an invariant camera space. Doing so enables two camera systems with nearly the same visual evidence of an object to transform into the same virtual depth as shown in Figure 4 (Main paper). Next we provide a simple proof for the conversion between virtual and metric depth.

*Proof:* Assume a 3D point  $(X, Y, Z)$  projected to  $(x, y)$  on an image with height  $H$  and focal length  $f$ . The virtual 3D point  $(X, Y, Z_v)$  is projected to  $(x_v, y_v)$  on the virtual image. The 2D points  $(x, y)$  and  $(x_v, y_v)$  correspond to the same pixel location in both the original and virtual image. In other words,  $y_v = y \cdot \frac{H_v}{H}$ . Recall the formula for projection as  $y = f \cdot \frac{Y}{Z} + p_y$  and  $y_v = f_v \cdot \frac{Y_v}{Z_v} + p_{y_v}$ , where  $p_y$  is the principal point and  $p_{y_v} = p_y \cdot \frac{H_v}{H}$ . By substitution  $f_v \cdot \frac{Y_v}{Z_v} = f \cdot \frac{Y}{Z} \cdot \frac{H_v}{H} \Rightarrow Z_v = Z \cdot \frac{f_v}{f} \cdot \frac{H}{H_v}$ .

## A2.3 Model Efficiency

When training on subsets smaller than OMNI3D, we adjust the learning rate, batch size, and number of iterations linearly until we can train for 128 epochs between 96k to 116k iterations. Cube R-CNN trains on V100 GPUs between 14 and 26 hours depending on the subset configuration when scaled to multi-node distributed training, and while training uses approximately 1.6 GB memory per image. *Inference* on a Cube R-CNN model processes image from KITTI [19] with a wide input resolution of  $512 \times 1696$  at 52ms/image on average while taking up 1.3 GB memory on a Quadro GP100 GPU. Computed with an identical environment, our model efficiency is favorable to M3D-RPN [6] and GUPNet [49] which infer from KITTI images at 191ms and 66ms on average, respectively.

## A3 Evaluation

In this section we give additional context and justification for our chosen  $\tau$  which defines the range of 3D IoU that  $\text{AP}_{3D}$  is averaged over. We further provide details on our implementation of 3D IoU.

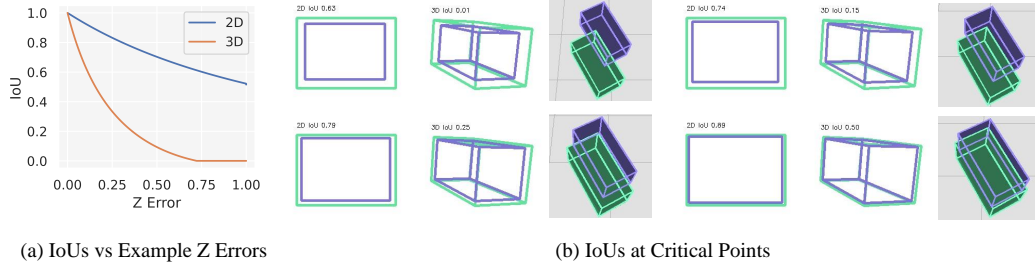


Figure 9: We slowly translate a 3D box (purple) backwards relative to its ground-truth (green), hence simulating error in  $z$  up to 1 unit length ( $\bar{l}_{3D}^{gt} = 1$ ). We plot 2D and 3D IoU vs  $z$  error in (a) and visualize selected critical points showing their 2D projections, 3D front view and a novel top view. Unsurprisingly, we find the drop in  $\text{IoU}_{3D}$  is much steeper than 2D. This effect highlights both the challenge of 3D object detection and helps justify the relaxed  $\tau$  thresholds used for  $\text{AP}_{3D}$  in Sec. 5.

**Data:** Two 3D boxes  $b_1$  and  $b_2$

**Result:** Intersecting shape  $S = []$

Step 1: For each 3D triangular face  $e \in b_1$  we check whether  $e$  falls inside  $b_2$

Step 2: If  $e$  is not inside, then we discard it

Step 3: If  $e$  is inside, then  $S = S + [e]$ . If  $e$  is partially inside, then the part of  $e$  inside  $b_2$ , call it  $\hat{e}$ , is added to  $S$ ,  $S = S + [\hat{e}]$

Step 4: We repeat steps 1 - 3 for  $b_2$

Step 5: We check and remove duplicates in  $S$  (in the case of coplanar sides in  $b_1$  and  $b_2$ )

Step 6: We compute the volume of  $S$ , which is guaranteed to be convex

**Algorithm 1:** An high-level overview of our fast and exact  $\text{IoU}_{3D}$  algorithm.

### A3.1 Thresholds $\tau$ for $\text{AP}_{3D}$

We highlight the relationship between 2D and 3D IoU in Figure 9. We do so by contriving a general example between a ground-truth and a predicted box, both of which share rotation and rectangular cuboid dimensions ( $0.5 \times 0.5 \times 1.0$ ). We then translate the 3D box along the  $z$ -axis up to 1 meter (its unit length), simulating small to large errors in depth estimation. As shown in the left of Figure 9, the 3D IoU drops off significantly quicker than 2D IoU does between the projected 2D bounding boxes. As visualized in right of Figure 9, a moderate score of 0.63  $\text{IoU}_{2D}$  may result in a low 0.01  $\text{IoU}_{3D}$ . Despite visually appearing to be well localized in the front view, the top view helps reveal the error. Since depth is a key error source for 3D, we find the relaxed settings of  $\tau$  (Sec. 5) to be reasonable.

### A3.2 $\text{IoU}_{3D}$ Details

We implement a fast  $\text{IoU}_{3D}$  implementation. We provide more details for our algorithm here. Our algorithm starts from the simple observation that the intersection of two oriented 3D boxes,  $b_1$  and  $b_2$ , is a convex polyhedron with  $n > 2$  comprised of connected planar units. In 3D, these planar units are 3D triangular faces. Critically, each planar unit belongs strictly to either  $b_1$  or  $b_2$ . Our algorithm finds these units by iterating through the sides of each box, as described in Algorithm 1.

## A4 Qualitative Examples

We provide more examples of 3D object detection from Cube R-CNN in Figure 10 from OMNI3D test. In Figure 11, we demonstrate generalization for interesting scenes in the wild from COCO [42] images. When projecting on images with unknown camera intrinsics, as is the case for COCO, we visualize with intrinsics of  $f = 2 \cdot H$ ,  $p_x = \frac{1}{2}W$ ,  $p_y = \frac{1}{2}H$ , where  $H \times W$  is the input image resolution. As shown in Fig. 11, this appears to result in fairly stable generalization for indoor and more common failure cases concerning unseen object or camera poses in outdoor.

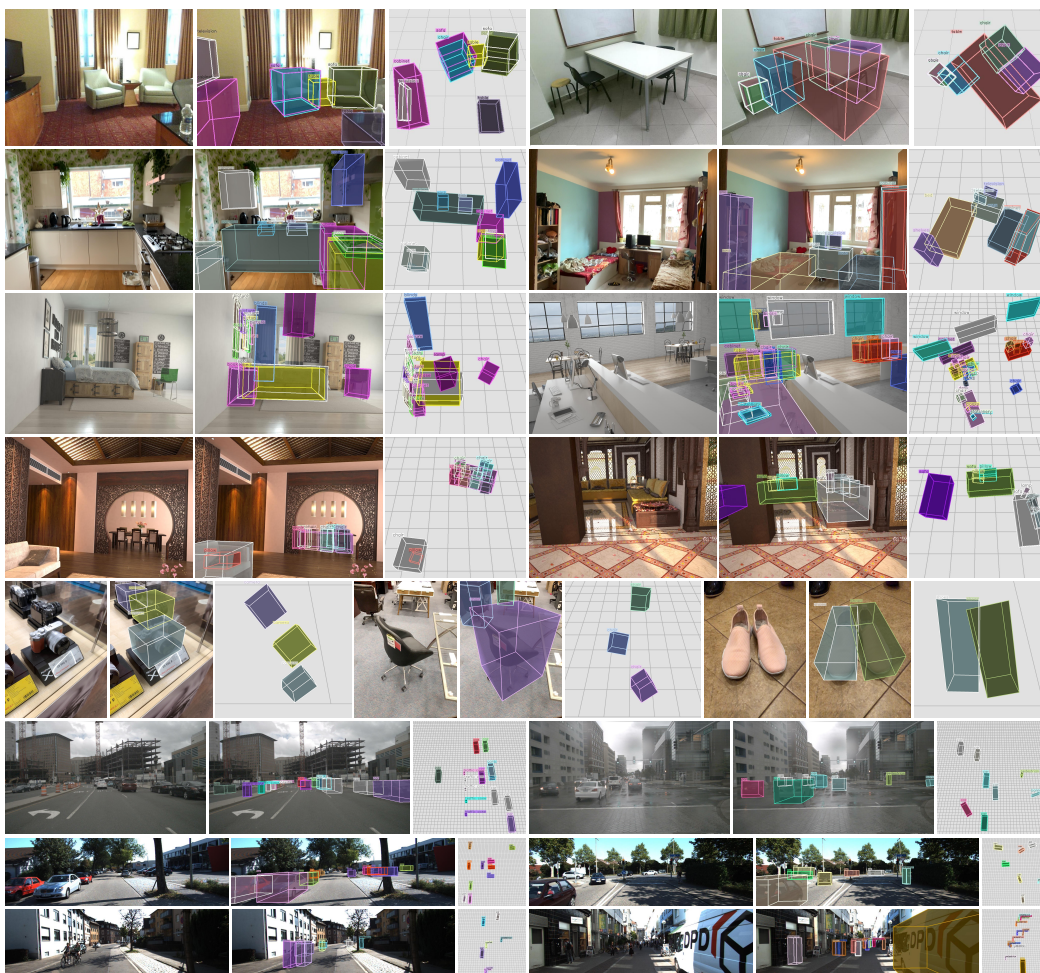


Figure 10: Cube R-CNN on OMN13D test. We show the input image, the 3D predictions overlaid on the image and a top view. Rows contains samples from datasets in order of SUN RGB-D [66], ARKit [5], Hypersim [62] (Rows 3–4), Objectron [2], nuScenes [8], and KITTI [19] (Rows 7–8).



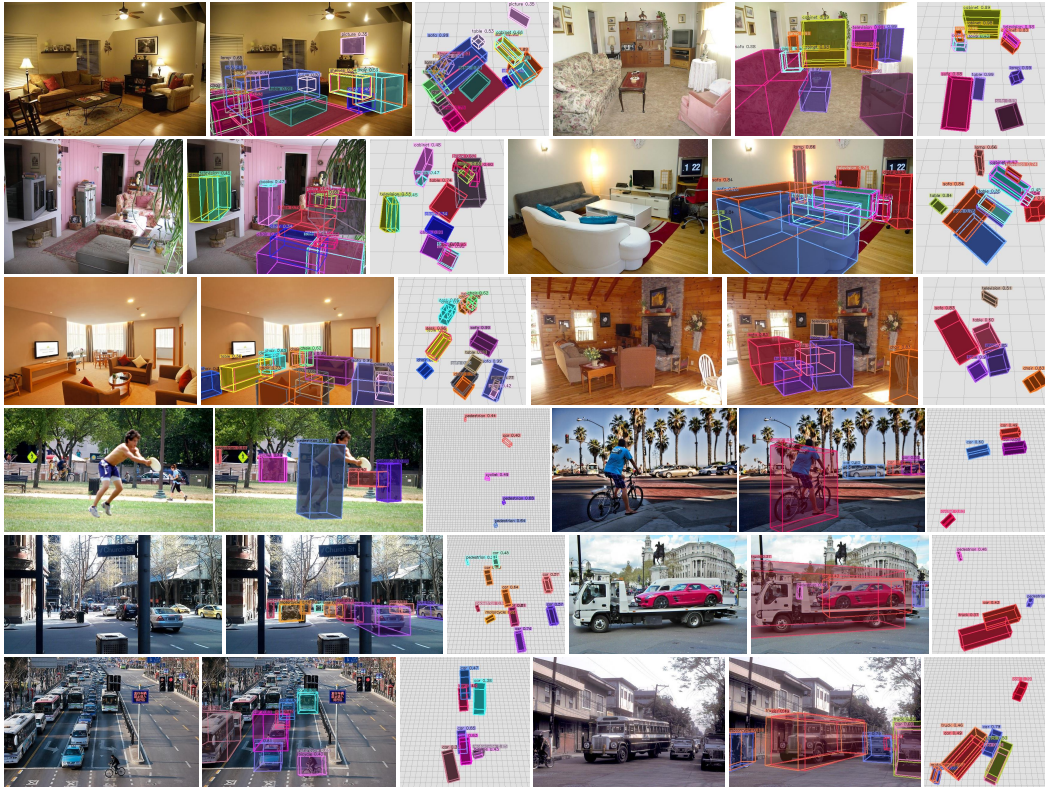


Figure 11: Cube R-CNN on COCO in the wild images. We select interesting scenes and observe that generalization tends to perform well for indoor scenes (Rows 1-3) compared to outdoor (Rows 4-6), which appear to fail when in the wild objects or cameras have high variation in unseen poses.



## References

- [1] Waymo open dataset: An autonomous driving dataset (2019)
- [2] Ahmadyan, A., Zhang, L., Ablavatski, A., Wei, J., Grundmann, M.: Objectron: A large scale dataset of object-centric videos in the wild with pose annotations. CVPR (2021)
- [3] Bao, W., Xu, B., Chen, Z.: MonoFENet: Monocular 3D object detection with feature enhancement networks. TIP **29**, 2753–2765 (2019)
- [4] Barber, C.B., Dobkin, D.P., Huhdanpaa, H.: The quickhull algorithm for convex hulls. ACM Transactions on Mathematical Software (TOMS) **22**(4), 469–483 (1996)
- [5] Baruch, G., Chen, Z., Dehghan, A., Dimry, T., Feigin, Y., Fu, P., Gebauer, T., Joffe, B., Kurz, D., Schwartz, A., Shulman, E.: ARKitscenes - a diverse real-world dataset for 3d indoor scene understanding using mobile RGB-d data. In: NeurIPS Datasets and Benchmarks Track (Round 1) (2021)
- [6] Brazil, G., Liu, X.: M3D-RPN: Monocular 3D region proposal network for object detection. In: ICCV (2019)
- [7] Brazil, G., Pons-Moll, G., Liu, X., Schiele, B.: Kinematic 3D object detection in monocular video. In: ECCV. pp. 135–152. Springer (2020)
- [8] Caesar, H., Bankiti, V., Lang, A.H., Vora, S., Liong, V.E., Xu, Q., Krishnan, A., Pan, Y., Baldan, G., Beijbom, O.: nuscenes: A multimodal dataset for autonomous driving. In: CVPR (2020)
- [9] Chang, M.F., Lambert, J., Sangkloy, P., Singh, J., Bak, S., Hartnett, A., Wang, D., Carr, P., Lucey, S., Ramanan, D., et al.: Argoverse: 3d tracking and forecasting with rich maps. In: CVPR. pp. 8748–8757 (2019)
- [10] Chen, H., Huang, Y., Tian, W., Gao, Z., Xiong, L.: Monorun: Monocular 3D object detection by reconstruction and uncertainty propagation. In: CVPR. pp. 10379–10388 (2021)
- [11] Chen, K., Wang, J., Pang, J., Cao, Y., Xiong, Y., Li, X., Sun, S., Feng, W., Liu, Z., Xu, J., Zhang, Z., Cheng, D., Zhu, C., Cheng, T., Zhao, Q., Li, B., Lu, X., Zhu, R., Wu, Y., Dai, J., Wang, J., Shi, J., Ouyang, W., Loy, C.C., Lin, D.: MMDetection: Open mmlab detection toolbox and benchmark. arXiv preprint arXiv:1906.07155 (2019)
- [12] Chen, X., Kundu, K., Zhang, Z., Ma, H., Fidler, S., Urtasun, R.: Monocular 3D object detection for autonomous driving. In: CVPR (2016)
- [13] Chen, Y.N., Dai, H., Ding, Y.: Pseudo-stereo for monocular 3D object detection in autonomous driving. In: CVPR (2022)
- [14] Chen, Y., Tai, L., Sun, K., Li, M.: Monopair: Monocular 3D object detection using pairwise spatial relationships. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12093–12102 (2020)
- [15] Dasgupta, S., Fang, K., Chen, K., Savarese, S.: Delay: Robust spatial layout estimation for cluttered indoor scenes. In: CVPR (2016)
- [16] Ding, M., Huo, Y., Yi, H., Wang, Z., Shi, J., Lu, Z., Luo, P.: Learning depth-guided convolutions for monocular 3D object detection. In: CVPRW. pp. 1000–1001 (2020)
- [17] Fu, H., Gong, M., Wang, C., Batmanghelich, K., Tao, D.: Deep Ordinal Regression Network for Monocular Depth Estimation. In: CVPR (2018)
- [18] Gählert, N., Jourdan, N., Cordts, M., Franke, U., Denzler, J.: Cityscapes 3D: Dataset and benchmark for 9 DOF vehicle detection. In: CVPRW (2020)
- [19] Geiger, A., Lenz, P., Urtasun, R.: Are we ready for autonomous driving? the kitti vision benchmark suite. In: CVPR (2012)
- [20] Gkioxari, G., Malik, J., Johnson, J.: Mesh R-CNN. In: ICCV (2019)

- [21] Gu, J., Wu, B., Fan, L., Huang, J., Cao, S., Xiang, Z., Hua, X.S.: Homography loss for monocular 3d object detection. In: CVPR (2022)
- [22] Gupta, A., Dollar, P., Girshick, R.: LVIS: A dataset for large vocabulary instance segmentation. In: CVPR (2019)
- [23] He, C., Huang, J., Hua, X.S., Zhang, L.: Aug3d-rpn: Improving monocular 3d object detection by synthetic images with virtual depth. arXiv preprint arXiv:2107.13269 (2021)
- [24] He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask R-CNN. In: ICCV (2017)
- [25] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016)
- [26] Hedau, V., Hoiem, D., Forsyth, D.: Recovering the spatial layout of cluttered rooms. In: ICCV (2009)
- [27] Houston, J., Zuidhof, G., Bergamini, L., Ye, Y., Chen, L., Jain, A., Omari, S., Iglovikov, V., Ondruska, P.: One thousand and one hours: Self-driving motion prediction dataset. arXiv preprint arXiv:2006.14480 (2020)
- [28] Huang, K.C., Wu, T.H., Su, H.T., Hsu, W.H.: Monodtr: Monocular 3D object detection with depth-aware transformer. In: CVPR (2022)
- [29] Huang, S., Qi, S., Xiao, Y., Zhu, Y., Wu, Y.N., Zhu, S.C.: Cooperative holistic scene understanding: Unifying 3D object, layout, and camera pose estimation. In: NeurIPS (2018)
- [30] Huang, X., Wang, P., Cheng, X., Zhou, D., Geng, Q., Yang, R.: The apolloscape open dataset for autonomous driving and its application. PAMI **42**(10), 2702–2719 (2019)
- [31] Kesten, R., Usman, M., Houston, J., Pandya, T., Nadhamuni, K., Ferreira, A., Yuan, M., Low, B., Jain, A., Ondruska, P., Omari, S., Shah, S., Kulkarni, A., Kazakova, A., Tao, C., Platinsky, L., Jiang, W., Shet, V.: Level 5 perception dataset 2020. <https://level-5.global/level5/data/> (2019)
- [32] Kim, D., Lin, T.Y., Angelova, A., Kweon, I.S., Kuo, W.: Learning open-world object proposals without learning to classify. IEEE Robotics and Automation Letters **7**(2), 5453–5460 (2022)
- [33] Ku, J., Mozifian, M., Lee, J., Harakeh, A., Waslander, S.L.: Joint 3D proposal generation and object detection from view aggregation. In: IROS. pp. 1–8. IEEE (2018)
- [34] Kulkarni, N., Misra, I., Tulsiani, S., Gupta, A.: 3d-relnet: Joint object and relational network for 3d prediction. In: ICCV (2019)
- [35] Kumar, A., Brazil, G., Liu, X.: Groomed-nms: Grouped mathematically differentiable nms for monocular 3D object detection. In: CVPR. pp. 8973–8983 (2021)
- [36] Kundu, A., Li, Y., Rehg, J.M.: 3d-rcnn: Instance-level 3d object reconstruction via render-and-compare. In: CVPR. pp. 3559–3568 (2018)
- [37] Law, H., Deng, J.: Cornernet: Detecting objects as paired keypoints. In: ECCV. pp. 734–750 (2018)
- [38] Lee, D.C., Hebert, M., Kanade, T.: Geometric reasoning for single image structure recovery. In: CVPR (2009)
- [39] Li, P., Zhao, H., Liu, P., Cao, F.: RTM3D: Real-time monocular 3D detection from object keypoints for autonomous driving. In: ECCV. pp. 644–660. Springer (2020)
- [40] Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: CVPR (2017)
- [41] Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: ICCV (2017)

- [42] Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft COCO: Common objects in context. In: ECCV (2014)
- [43] Liu, C., Gu, S., Van Gool, L., Timofte, R.: Deep line encoding for monocular 3D object detection and depth prediction. In: BMVC 2021. p. 354 (2021)
- [44] Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C.: Ssd: Single shot multibox detector. In: ECCV (2016)
- [45] Liu, X., Xue, N., Wu, T.: Learning auxiliary monocular contexts helps monocular 3d object detection. In: AAAI (2022)
- [46] Liu, Y., Yixuan, Y., Liu, M.: Ground-aware monocular D object detection for autonomous driving. RA-L 6(2), 919–926 (2021)
- [47] Liu, Z., Wu, Z., Tóth, R.: Smoke: Single-stage monocular 3D object detection via keypoint estimation. In: CVPR W (2020)
- [48] Liu, Z., Zhou, D., Lu, F., Fang, J., Zhang, L.: Autoshape: Real-time shape-aware monocular 3D object detection. In: ICCV. pp. 15641–15650 (2021)
- [49] Lu, Y., Ma, X., Yang, L., Zhang, T., Liu, Y., Chu, Q., Yan, J., Ouyang, W.: Geometry uncertainty projection network for monocular 3D object detection. In: ICCV. pp. 3111–3121 (2021)
- [50] Ma, X., Liu, S., Xia, Z., Zhang, H., Zeng, X., Ouyang, W.: Rethinking pseudo-lidar representation. In: ECCV (2020)
- [51] Ma, X., Wang, Z., Li, H., Zhang, P., Ouyang, W., Fan, X.: Accurate monocular 3D object detection via color-embedded 3D reconstruction for autonomous driving. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 6851–6860 (2019)
- [52] Ma, X., Zhang, Y., Xu, D., Zhou, D., Yi, S., Li, H., Ouyang, W.: Delving into localization errors for monocular 3D object detection. In: CVPR (June 2021)
- [53] Mallya, A., Lazebnik, S.: Learning informative edge maps for indoor scene layout prediction. In: ICCV (2015)
- [54] Mousavian, A., Anguelov, D., Flynn, J., Kosecka, J.: 3D bounding box estimation using deep learning and geometry. In: CVPR (2017)
- [55] Nie, Y., Han, X., Guo, S., Zheng, Y., Chang, J., Zhang, J.J.: Total3DUnderstanding: Joint layout, object pose and mesh reconstruction for indoor scenes from a single image. In: CVPR (2020)
- [56] Park, D., Ambrus, R., Guizilini, V., Li, J., Gaidon, A.: Is pseudo-lidar needed for monocular 3D object detection? In: ICCV. pp. 3142–3152 (2021)
- [57] Qi, C.R., Liu, W., Wu, C., Su, H., Guibas, L.J.: Frustum pointnets for 3D object detection from rgb-d data. In: CVPRn (2018)
- [58] Ravi, N., Reizenstein, J., Novotny, D., Gordon, T., Lo, W.Y., Johnson, J., Gkioxari, G.: Accelerating 3D deep learning with PyTorch3D. arXiv:2007.08501 (2020)
- [59] Reading, C., Harakeh, A., Chae, J., Waslander, S.L.: Categorical depth distribution network for monocular 3D object detection. In: CVPR. pp. 8555–8564 (2021)
- [60] Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: Unified, real-time object detection. In: CVPR (2016)
- [61] Ren, S., He, K., Girshick, R., Sun, J.: Faster R-CNN: Towards real-time object detection with region proposal networks. In: NeurIPS (2015)
- [62] Roberts, M., Ramapuram, J., Ranjan, A., Kumar, A., Bautista, M.A., Paczan, N., Webb, R., Susskind, J.M.: Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding. In: ICCV (2021)

- [63] Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. *IJCV* (2015)
- [64] Simonelli, A., Bulò, S.R., Porzi, L., López-Antequera, M., Kotschieder, P.: Disentangling monocular 3D object detection. In: *ICCV* (2019)
- [65] Simonelli, A., Bulò, S.R., Porzi, L., Ricci, E., Kotschieder, P.: Towards generalization across depth for monocular 3D object detection. In: *ECCV* (2020)
- [66] Song, S., Lichtenberg, S.P., Xiao, J.: Sun rgb-d: A rgb-d scene understanding benchmark suite. In: *CVPR* (2015)
- [67] Tian, Z., Shen, C., Chen, H., He, T.: FCOS: Fully convolutional one-stage object detection. In: *ICCV*. pp. 9627–9636 (2019)
- [68] Tulsiani, S., Gupta, S., Fouhey, D., Malik, A.A.E.J.: Factoring shape, pose, and layout from the 2d image of a 3d scene. In: *CVPR* (2018)
- [69] Virtanen, P., Gommers, R., Oliphant, T.E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., et al.: Scipy 1.0: fundamental algorithms for scientific computing in python. *Nature methods* **17**(3), 261–272 (2020)
- [70] Wang, L., Du, L., Ye, X., Fu, Y., Guo, G., Xue, X., Feng, J., Zhang, L.: Depth-conditioned dynamic message propagation for monocular 3d object detection. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 454–463 (2021)
- [71] Wang, L., Zhang, L., Zhu, Y., Zhang, Z., He, T., Li, M., Xue, X.: Progressive coordinate transforms for monocular 3D object detection. In: *NeurIPS*. vol. 34 (2021)
- [72] Wang, Y., Chao, W.L., Garg, D., Hariharan, B., Campbell, M., Weinberger, K.Q.: Pseudo-lidar from visual depth estimation: Bridging the gap in 3D object detection for autonomous driving. In: *CVPR*. pp. 8445–8453 (2019)
- [73] Wu, Y., Kirillov, A., Massa, F., Lo, W.Y., Girshick, R.: Detectron2. <https://github.com/facebookresearch/detectron2> (2019)
- [74] You, Y., Wang, Y., Chao, W.L., Garg, D., Pleiss, G., Hariharan, B., Campbell, M., Weinberger, K.Q.: Pseudo-lidar++: Accurate depth for 3D object detection in autonomous driving. In: *ICLR* (2020)
- [75] Yu, F., Wang, D., Shelhamer, E., Darrell, T.: Deep layer aggregation. In: *CVPR*. pp. 2403–2412 (2018)
- [76] Zhang, Y., Lu, J., Zhou, J.: Objects are different: Flexible monocular 3D object detection. In: *CVPR*. pp. 3289–3298 (2021)
- [77] Zhou, D., Song, X., Dai, Y., Yin, J., Lu, F., Liao, M., Fang, J., Zhang, L.: Iafa: Instance-aware feature aggregation for 3D object detection from a single image. In: *ACCV* (2020)
- [78] Zhou, X., Wang, D., Krähenbühl, P.: Objects as points. In: *arXiv preprint arXiv:1904.07850* (2019)
- [79] Zhou, Y., Barnes, C., Lu, J., Yang, J., Li, H.: On the continuity of rotation representations in neural networks. In: *CVPR* (2019)
- [80] Zhou, Y., He, Y., Zhu, H., Wang, C., Li, H., Jiang, Q.: Monocular 3D object detection: An extrinsic parameter free approach. In: *CVPR*. pp. 7556–7566 (2021)
- [81] Zou, Z., Ye, X., Du, L., Cheng, X., Tan, X., Zhang, L., Feng, J., Xue, X., Ding, E.: The devil is in the task: Exploiting reciprocal appearance-localization features for monocular 3D object detection. In: *ICCV*. pp. 2713–2722 (2021)