# DEVIANT: Depth EquiVarIAnt NeTwork for Monocular 3D Object Detection

Abhinav Kumar[1], Garrick Brazil[2], Enrique Corona[3], Armin Parchami[3], and Xiaoming Liu[1]

[1] Michigan State University
[2] Meta AI
[3] Ford Motor Company

[1]{kumarab6, liuxm}@msu.edu,[2]brazilga@fb.com,[3]{ecoron18, mparcham}@ford.com
https://github.com/abhi1kumar/DEVIANT

**Abstract.** Modern neural networks use building blocks such as convolutions that are equivariant to arbitrary 2D translations. However, these vanilla blocks are not equivariant to arbitrary 3D translations in the projective manifold. Even then, all monocular 3D detectors use vanilla blocks to obtain the 3D coordinates, a task for which the vanilla blocks are not designed for. This paper takes the first step towards convolutions equivariant to arbitrary 3D translations in the projective manifold. Since the depth is the hardest to estimate for monocular detection, this paper proposes Depth EquiVarIAnt NeTwork (DEVIANT) built with existing scale equivariant steerable blocks. As a result, DEVIANT is equivariant to the depth translations in the projective manifold whereas vanilla networks are not. The additional depth equivariance forces the DEVIANT to learn consistent depth estimates, and therefore, DEVIANT achieves state-of-the-art monocular 3D detection results on KITTI and Waymo datasets in the image-only category and performs competitively to methods using extra information. Moreover, DEVIANT works better than vanilla networks in cross-dataset evaluation.

**Keywords:** Equivariance, Projective manifold, Monocular 3D detection

## 1 Introduction

Monocular 3D object detection is a fundamental task in computer vision, where the task is to infer 3D information including depth from a single monocular image. It has applications in augmented reality [2], gaming [63], robotics [65], and more recently in autonomous driving [4,68] as a fallback solution for LiDAR.

Most of the monocular 3D methods attach extra heads to the 2D Faster-RCNN [64] or CenterNet [102] for 3D detections. Some change architectures [42,45,76] or losses [4,13]. Others incorporate augmentation [71], or confidence [5,45]. Recent ones use in-network ensembles [49,100] for better depth estimation.

Most of these methods use vanilla blocks such as convolutions that are *equivariant* to arbitrary 2D translations [6,61]. In other words, whenever we shift the ego camera in 2D (See $t_u$ of Fig. 1), the new image (projection) is a translation of
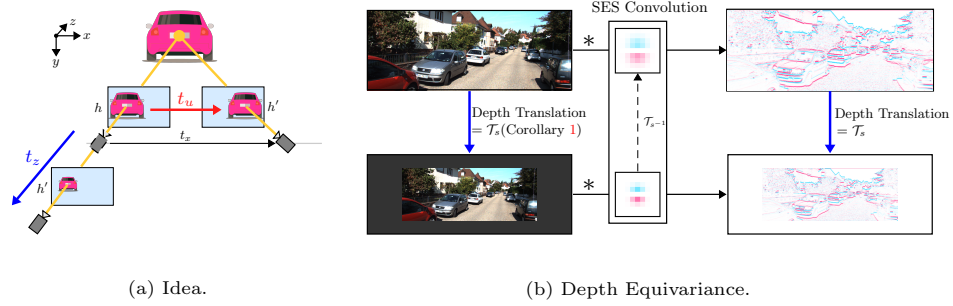
(a) Idea.                                    (b) Depth Equivariance.

Fig. 1: **(a) Idea.** Vanilla CNN is equivariant to projected 2D translations $t_u, t_v$ of the ego camera. The ego camera moves in 3D in driving scenes which breaks this assumption. We propose DEVIANT which is additionally equivariant to depth translations $t_z$ in the projective manifold. **(b) Depth Equivariance**. DEVIANT enforces additional consistency among the feature maps of an image and its transformation caused by the ego depth translation. $\mathcal{T}_s$ = scale transformation, $*$ = vanilla convolution.

the original image, and therefore, these methods output a translated feature map. However, in general, the camera moves in depth in driving scenes instead of 2D (See $t_z$ of Fig. 1). So, the new image is not a translation of the original input image due to the projective transform. Thus, using vanilla blocks in monocular methods is a mismatch between the assumptions

Table 1: **Equivariance comparisons**. [Key: Proj.= Projected, ax= axis]

| **Translation** ➜ | 3D | | | Proj. 2D | |
|---|---|---|---|---|---|
| | $x$−ax | $y$−ax | $z$−ax | $u$-ax | $v$-ax |
| | $(t_x)$ | $(t_y)$ | $(t_z)$ | $(t_u)$ | $(t_v)$ |
| Vanilla CNN | – | – | – | ✓ | ✓ |
| Log-polar [106] | – | – | ✓ | – | – |
| **DEVIANT** | – | – | ✓ | ✓ | ✓ |
| Ideal | ✓ | ✓ | ✓ | – | – |

and the regime where these blocks operate. Additionally, there is a huge generalization gap between training and validation for monocular 3D detection (See Tab. 14 in the supplementary). Modeling translation equivariance in the correct manifold improves generalization for tasks in spherical [15] and hyperbolic [26] manifolds. Monocular detection involves processing pixels (3D point projections) to obtain the 3D information, and is thus a task in the projective manifold. Moreover, the depth in monocular detection is ill-defined [76], and thus, the hardest to estimate [53]. Hence, using building blocks *equivariant to depth translations in the projective manifold* is a natural choice for improving generalization and is also at the core of this work (See Appendix A1.8).

Recent monocular methods use flips [4], scale [49,71], mosaic [3,77] or copy-paste [43] augmentation, depth-aware convolution [4], or geometry [47,49,67,99] to improve generalization. Although all these methods improve performance, a major issue is that their backbones are not designed for the projective world. This results in the depth estimation going haywire with a slight ego movement [103]. Moreover, data augmentation, *e.g.*, flips, scales, mosaic, copy-paste, is not only limited for the projective tasks, but also does not guarantee desired behavior [25].

To address the mismatch between assumptions and the operating regime of the vanilla blocks and improve generalization, we take the first step towards convolutions equivariant to arbitrary 3D translations in the projective mani-

Table 2: Equivariances known in the literature.

| Transformation → <br> Manifold ↓ | Translation | Rotation | Scale | Flips | Learned |
|---|---|---|---|---|---|
| Euclidean | Vanilla CNN [40] | Polar, Steerable [91] | Log-polar [31], Steerable [29] | ChiralNets [96] | Transformers [21] |
| Spherical | Spherical CNN [15] | – | – | – | – |
| Hyperbolic | Hyperbolic CNN [26] | – | – | – | – |
| Projective | Monocular Detector | – | – | – | – |

fold. We propose Depth EquiVarIAnt NeTwork (DEVIANT) which is additionally equivariant to depth translations in the projective manifold as shown in Tab. 1. Building upon the classic result from [30], we simplify it under reasonable assumptions about the camera movement in autonomous driving to get scale transformations. The scale equivariant blocks are well-known in the literature [29, 32, 74, 104], and consequently, we replace the vanilla blocks in the backbone with their scale equivariant steerable counterparts [74] to additionally embed equivariance to depth translations in the projective manifold. Hence, DEVIANT learns consistent depth estimates and improves monocular detection.

In summary, the main contributions of this work include:

- We study the modeling error in monocular 3D detection and propose depth equivariant networks built with scale equivariant steerable blocks as a solution.
- We achieve state-of-the-art (SOTA) monocular 3D object detection results on the KITTI and Waymo datasets in the image-only category and perform competitively to methods which use extra information.
- We experimentally show that DEVIANT works better in cross-dataset evaluation suggesting better generalization than vanilla CNN backbones.

## 2 Literature Review

**Equivariant Neural Networks.** The success of convolutions in CNN has led people to look for their generalizations [17, 87]. Convolution is the unique solution to 2D translation equivariance in the Euclidean manifold [6, 7, 61]. Thus, convolution in CNN is a prior in the Euclidean manifold. Several works explore other group actions in the Euclidean manifold such as 2D rotations [16, 19, 55, 88], scale [34, 54], flips [96], or their combinations [81, 91]. Some consider 3D translations [90] and rotations [78]. Few [21, 89, 101] attempt learning the equivariance from the data, but such methods have significantly higher data requirements [90]. Others change the manifold to spherical [15], hyperbolic [26], graphs [56], or arbitrary manifolds [33]. Monocular 3D detection involves operations on pixels which are projections of 3D point and thus, works in a different manifold namely projective manifold. Tab. 2 summarizes all these equivariances known thus far.

**Scale Equivariant Networks.** Scale equivariance in the Euclidean manifold is more challenging than the rotations because of its acyclic and unbounded nature [61]. There are two major lines of work for scale equivariant networks. The first [22, 31] infers the global scale using log-polar transform [106], while the other infers the scale locally by convolving with multiple scales of images [34] or

filters [94]. Several works [29, 32, 74, 104] extend the local idea, using steerable filters [24]. Another work [92] constructs filters for integer scaling. We compare the two kinds of scale equivariant convolutions on the monocular 3D detection task and show that steerable convolutions are better suited to embed depth (scale) equivariance. Scale equivariant networks have been used for classification [22, 29, 74], 2D tracking [73] and 3D object classification [22]. We are the first to use scale equivariant networks for monocular 3D detection.

**3D Object Detection.** Accurate 3D object detection uses sparse data from LiDARs [66], which are expensive and do not work well in severe weather [76] and glassy environments. Hence, several works have been on monocular camera-based 3D object detection, which is simplistic but has scale/depth ambiguity [76]. Earlier approaches [11, 23, 59, 60] use hand-crafted features, while the recent ones use deep learning. Some change architectures [42, 45, 46, 76] or losses [4, 13]. Some use scale [49, 71], mosaic [77] or copy-paste [43] augmentation. Others incorporate depth in convolution [4, 20], or confidence [5, 37, 45]. More recent ones use in-network ensembles to predict the depth deterministically [100] or probabilistically [49]. A few use temporal cues [5], NMS [36], or corrected camera extrinsics [103] in the training pipeline. Some also use CAD models [10, 48] or LiDAR [62] in training. Another line of work called Pseudo-LiDAR [50, 52, 57, 69, 83] estimates the depth first, and then uses a point cloud-based 3D object detector. We refer to [51] for a detailed survey. Our work is the first to use scale equivariant blocks in the backbone for monocular 3D detection.

## 3    Background

We first provide the necessary definitions which are used throughout this paper. These are not our contributions and can be found in the literature [8, 30, 90].

**Equivariance.** Consider a group of transformations $G$, whose individual members are $g$. Assume $\Phi$ denote the mapping of the inputs $h$ to the outputs $y$. Let the inputs and outputs undergo the transformation $\mathcal{T}_g^h$ and $\mathcal{T}_g^y$ respectively. Then, the mapping $\Phi$ is equivariant to the group $G$ [90] if $\Phi(\mathcal{T}_g^h h) = \mathcal{T}_g^y(\Phi h), \forall\ g \in G$. Thus, equivariance provides an explicit relationship between input transformations and feature-space transformations at each layer of the neural network [90], and intuitively makes the learning easier. The mapping $\Phi$ is the vanilla convolution when the $\mathcal{T}_g^h = \mathcal{T}_g^y = \mathcal{T}_\mathbf{t}$ where $\mathcal{T}_\mathbf{t}$ denotes the translation $\mathbf{t}$ on the discrete grid [6, 7, 61]. These vanilla convolution introduce weight-tying [40] in fully connected neural networks resulting in a greater generalization. A special case of equivariance is the invariance [90] which is given by $\Phi(\mathcal{T}_g^h h) = \Phi h, \forall\ g \in G$.

**Projective Transformations.** Our idea is to use equivariance to depth translations in the projective manifold since the monocular detection task belongs to this manifold. A natural question to ask is whether such equivariants exist in the projective manifold. [8] answers this question in negative, and says that such equivariants do not exist in general. However, such equivariants exist for special classes, such as planes. An intuitive way to understand this is to infer the rotations and translations by looking at the two projections (images). For

example, the result of [8] makes sense if we consider a car with very different front and back sides as in Fig. 6. A 180° ego rotation around the car means the projections (images) are its front and the back sides, which are different. Thus, we can not infer the translations and rotations from these two projections. Based on this result, we stick with **locally** planar objects *i.e.* we assume that a 3D object is made of several *patch planes*. (See last row of Fig. 2b as an example). It is important to stress that we do **NOT** assume that the 3D object such as car is planar. The local planarity also agrees with the property of manifolds that manifolds locally resemble $n$-dimensional Euclidean space and because the projective transform maps planes to planes, the patch planes in 3D are also locally planar. We show a sample planar patch and the 3D object in Fig. 5 in the appendix.

**Planarity and Projective Transformation.** Example 13.2 from [30] links the planarity and projective transformations. Although their result is for stereo with two different cameras $(\mathbf{K}, \mathbf{K}')$, we substitute $\mathbf{K} = \mathbf{K}'$ to get Theorem 1.

**Theorem 1.** *[30] Consider a 3D point lying on a patch plane $mx+ny+oz+p=0$, and observed by an ego camera in a pinhole setup to give an image $h$. Let $\mathbf{t} = (t_x, t_y, t_z)$ and $\mathbf{R} = [r_{ij}]_{3 \times 3}$ denote a translation and rotation of the ego camera respectively. Observing the same 3D point from a new camera position leads to an image $h'$. Then, the image $h$ is related to the image $h'$ by the projective transformation*

$$\mathcal{T} : h(u - u_0, v - v_0) = \tag{1}$$

$$h'\left( f\frac{\left(r_{11}+\bar{t}_x\frac{m}{p}\right)(u-u_0)+\left(r_{21}+\bar{t}_x\frac{n}{p}\right)(v-v_0)+\left(r_{31}+\bar{t}_x\frac{o}{p}\right)f}{\left(r_{13}+\bar{t}_z\frac{m}{p}\right)(u-u_0)+\left(r_{23}+\bar{t}_z\frac{n}{p}\right)(v-v_0)+\left(r_{33}+\bar{t}_z\frac{o}{p}\right)f}, \right.$$

$$\left. f\frac{\left(r_{12}+\bar{t}_y\frac{m}{p}\right)(u-u_0)+\left(r_{22}+\bar{t}_y\frac{n}{p}\right)(v-v_0)+\left(r_{32}+\bar{t}_y\frac{o}{p}\right)f}{\left(r_{13}+\bar{t}_z\frac{m}{p}\right)(u-u_0)+\left(r_{23}+\bar{t}_z\frac{n}{p}\right)(v-v_0)+\left(r_{33}+\bar{t}_z\frac{o}{p}\right)f} \right),$$

*where $f$ and $(u_0, v_0)$ denote the focal length and principal point of the ego camera, and $(\bar{t}_x, \bar{t}_y, \bar{t}_z) = \mathbf{R}^T \mathbf{t}$.*

## 4    Depth Equivariant Backbone

The projective transformation in Eq. (1) from [30] is complicated and also involves rotations, and we do not know which convolution obeys this projective transformation. Hence, we simplify Eq. (1) under reasonable assumptions to obtain a familiar transformation for which the *convolution* is known.

**Corollary 1.** *When the ego camera translates in depth without rotations ($\mathbf{R} = \mathbf{I}$), and the patch plane is "approximately" parallel to the image plane, the image $h$ locally is a scaled version of the second image $h'$ independent of focal length, i.e.*

$$\mathcal{T}_s : h(u - u_0, v - v_0) \approx h'\left(\frac{u - u_0}{1 + t_z\frac{o}{p}}, \frac{v - v_0}{1 + t_z\frac{o}{p}}\right). \tag{2}$$

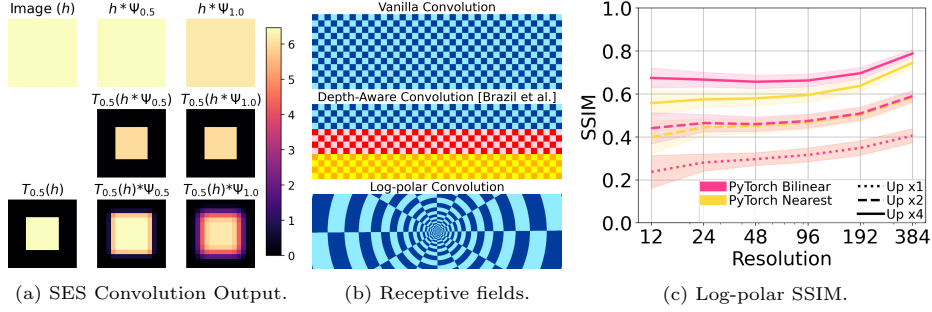(a) SES Convolution Output.    (b) Receptive fields.    (c) Log-polar SSIM.

Fig. 2: **(a) Scale Equivariance**. We apply SES convolution [74] with two scales on a single channel toy image $h$. **(b) Receptive fields** of convolutions in the Euclidean manifold. Colors represent different weights, while shades represent the same weight. **(c) Impact of discretization on log-polar convolution.** SSIM is very low at small resolutions and is not 1 even after upscaling by 4. [Key: Up= Upscaling]

*where $f$ and $(u_0, v_0)$ denote the focal length and principal point of the ego camera, and $t_z$ denotes the ego translation.*

See Appendix A1.6 for the detailed explanation of Corollary 1. Corollary 1 says

$$\mathcal{T}_s : h(u - u_0, v - v_0) \approx h'\left(\frac{u - u_0}{s}, \frac{v - v_0}{s}\right), \tag{3}$$

where, $s = 1 + t_z \frac{o}{p}$ denotes the scale and $\mathcal{T}_s$ denotes the scale transformation. The scale $s < 1$ suggests downscaling, while $s > 1$ suggests upscaling. Corollary 1 shows that the transformation $\mathcal{T}_s$ is independent of the focal length and that scale is a linear function of the depth translation. Hence, the depth translation in the projective manifold induces scale transformation and thus, the depth equivariance in the projective manifold is the scale equivariance in the Euclidean manifold. Mathematically, the desired equivariance is $[\mathcal{T}_s(h) * \Psi] = \mathcal{T}_s [h * \Psi_{s^{-1}}]$, where $\Psi$ denotes the filter (See Appendix A1.7). As CNN is not a scale equivariant (SE) architecture [74], we aim to get SE backbone which makes the architecture equivariant to depth translations in the projective manifold. The scale transformation is a familiar transformation and SE convolutions are well known [29, 32, 74, 104]. **Scale Equivariant Steerable (SES) Blocks.** We use the existing SES blocks [73, 74] to construct our Depth EquiVarIAnt NeTwork (DEVIANT) backbone. As [73] does not construct SE-DLA-34 backbones, we construct our DEVIANT backbone as follows. We replace the vanilla convolutions by the SES convolutions [73] with the basis as Hermite polynomials. SES convolutions result in multi-scale representation of an input tensor. As a result, their output is five-dimensional instead of four-dimensional. Thus, we replace the 2D pools and batch norm (BN) by 3D pools and 3D BN respectively. The Scale-Projection layer [74] carries a max over the extra (scale) dimension to project five-dimensional tensors to four dimensions (See Fig. 9 in the supplementary). Ablation in Sec. 5.2 confirms that BN and Pool (BNP) should also be SE for the best performance.

The SES convolutions [29, 74, 104] are based on steerable-filters [24]. Steerable approaches [29] first pre-calculate the non-trainable multi-scale basis in the Euclidean manifold and then build filters by the linear combinations of the trainable weights $\mathbf{w}$ (See Fig. 9). The number of trainable weights $\mathbf{w}$ equals the number of filters at one particular scale. The linear combination of multi-scale basis ensures that the filters are also multi-scale. Thus, SES blocks bypass grid conversion and do not suffer from sampling effects.

We show the convolution of toy image $h$ with a SES convolution in Fig. 2a. Let $\Psi_s$ denote the filter at scale $s$. The convolution between downscaled image and filter $\mathcal{T}_{0.5}(h) * \Psi_{0.5}$ matches the downscaled version of original image convolved with upscaled filter $\mathcal{T}_{0.5}(h * \Psi_{1.0})$. Fig. 2a (right column) shows that the output of a CNN exhibits aliasing in general and is therefore, not scale equivariant.

**Log-polar Convolution: Impact of Discretization.** An alternate way to convert the depth translation $t_z$ of Eq. (2) to shift is by converting the images to log-polar space [106] around the principal point $(u_0, v_0)$, as

$$h(\ln r, \theta) \approx h'\left(\ln r - \ln\left(1 + t_z \frac{o}{p}\right), \ \theta\right), \tag{4}$$

with $r = \sqrt{(u - u_0)^2 + (v - v_0)^2}$, and $\theta = \tan^{-1}\left(\frac{v - v_0}{u - u_0}\right)$. The log-polar transformation converts the scale to translation, so using convolution in the log-polar space is equivariant to the logarithm of the depth translation $t_z$. We show the receptive field of log-polar convolution in Fig. 2b. The log-polar convolution uses a smaller receptive field for objects closer to the principal point, while a larger field away from the principal point. We implemented log-polar convolution and found that its performance (See Tab. 11) is not acceptable, consistent with [74]. We attribute this behavior to the discretization of pixels and loss of 2D translation equivariance. Eq. (4) is perfectly valid in the continuous world (Note the use of parentheses instead of square brackets in Eq. (4)). However, pixels reside on discrete grids, which gives rise to sampling errors [38]. We discuss the impact of discretization on log-polar convolution in Sec. 5.2 and show it in Fig. 2c. Hence, we do not use log-polar convolution for the DEVIANT backbone.

**Comparison of Equivariances for Monocular 3D Detection.** We now compare equivariances for monocular 3D detection task. An ideal monocular detector should be equivariant to arbitrary 3D translations $(t_x, t_y, t_z)$. However, most monocular detectors [36, 49] estimate 2D projections of 3D centers and the depth, which they back-project in 3D world via known camera intrinsics. Thus, a good enough detector shall be equivariant to 2D translations $(t_u, t_v)$ for projected centers as well as equivariant to depth translations $(t_z)$.

Existing detector backbones [36, 49] are only equivariant to 2D translations as they use vanilla convolutions that produce 4D feature maps. Log-polar backbones is equivariant to logarithm of depth translations but not to 2D translations. DEVIANT uses SES convolutions to produce 5D feature maps. The extra dimension in 5D feature map captures the changes in scale (for depth), while these feature maps individually are equivariant to 2D translations (for projected centers). Hence, DEVIANT augments the 2D translation equivariance $(t_u, t_v)$

of the projected centers with the depth translation equivariance. We emphasize that although DEVIANT is **not** equivariant to arbitrary 3D translations in the projective manifold, DEVIANT **does** provide the equivariance to depth translations ($t_z$) and is thus a first step towards the ideal equivariance. Our experiments (Sec. 5) show that even this additional equivariance benefits monocular 3D detection task. This is expected because depth is the hardest parameter to estimate [53]. Tab. 1 summarizes these equivariances. Moreover, Tab. 10 empirically shows that 2D detection does not suffer and therefore, confirms that DEVIANT indeed augments the 2D equivariance with the depth equivariance. An idea similar to DEVIANT is the optical expansion [95] which augments optical flow with the scale information and benefits depth estimation.

## 5    Experiments

Our experiments use the KITTI [28], Waymo [75] and nuScenes datasets [9]. We modify the publicly-available PyTorch [58] code of GUP Net [49] and use the GUP Net model as our baseline. For DEVIANT, we keep the number of scales as three [73]. DEVIANT takes 8.5 hours for training and 0.04s per image for inference on a single A100 GPU. See Appendix A2.2 for more details.

**Evaluation Metrics.** KITTI evaluates on three object categories: Easy, Moderate and Hard. It assigns each object to a category based on its occlusion, truncation, and height in the image space. KITTI uses $AP_{3D|R_{40}}$ percentage metric on the Moderate category to benchmark models [28] following [68, 70].

Waymo evaluates on two object levels: Level_1 and Level_2. It assigns each object to a level based on the number of LiDAR points included in its 3D box. Waymo uses $APH_{3D}$ percentage metric which is the incorporation of heading information in $AP_{3D}$ to benchmark models. It also provides evaluation at three distances $[0, 30)$, $[30, 50)$ and $[50, \infty)$ meters.

**Data Splits.** We use the following splits of the KITTI,Waymo and nuScenes:
- *KITTI Test (Full) split*: Official KITTI 3D benchmark [1] consists of 7,481 training and 7,518 testing images [28].
- *KITTI Val split*: It partitions the 7,481 training images into 3,712 training and 3,769 validation images [12].
- *Waymo Val split*: This split [62,80] contains 52,386 training and 39,848 validation images from the front camera. We construct its training set by sampling every third frame from the training sequences as in [62,80].
- *nuScenes Val split:* It consists of 28,130 training and 6,019 validation images from the front camera [9]. We use this split for evaluation [67].

### 5.1    KITTI Test Monocular 3D Detection

**Cars.** Tab. 3 lists out the results of monocular 3D detection and BEV evaluation on KITTI Test cars. Tab. 3 results show that DEVIANT outperforms the GUP Net and several other SOTA methods on both tasks. Except DD3D [57] and

Table 3: **Results on KITTI Test cars** at $IoU_{3D} \geq 0.7$. Previous results are from the leader-board or papers. We show 3 methods in each Extra category and 6 methods in the image-only category. [Key: <span style="color:red">**Best**</span>, <span style="color:blue">**Second Best**</span>]

| Method | Extra | $AP_{3D|R_{40}}[\%](\blacktriangle)$ | | | $AP_{BEV|R_{40}}[\%](\blacktriangle)$ | | |
|---|---|---|---|---|---|---|---|
| | | Easy | Mod | Hard | Easy | Mod | Hard |
| AutoShape [48] | CAD | 22.47 | 14.17 | 11.36 | 30.66 | 20.08 | 15.59 |
| PCT [80] | Depth | 21.00 | 13.37 | 11.31 | 29.65 | 19.03 | 15.92 |
| DFR-Net [105] | Depth | 19.40 | 13.63 | 10.35 | 28.17 | 19.17 | 14.84 |
| MonoDistill [14] | Depth | 22.97 | 16.03 | 13.60 | 31.87 | 22.59 | 19.72 |
| PatchNet-C [69] | LiDAR | 22.40 | 12.53 | 10.60 | – | – | – |
| CaDDN [62] | LiDAR | 19.17 | 13.41 | 11.46 | 27.94 | 18.91 | 17.19 |
| DD3D [57] | LiDAR | 23.22 | 16.34 | 14.20 | 30.98 | 22.56 | 20.03 |
| MonoEF [103] | Odometry | 21.29 | 13.87 | 11.71 | 29.03 | 19.70 | 17.26 |
| Kinematic [5] | Video | 19.07 | 12.72 | 9.17 | 26.69 | 17.52 | 13.10 |
| GrooMeD-NMS [36] | – | 18.10 | 12.32 | 9.65 | 26.19 | 18.27 | 14.05 |
| MonoRCNN [67] | – | 18.36 | 12.65 | 10.03 | 25.48 | 18.11 | 14.10 |
| MonoDIS-M [68] | – | 16.54 | 12.97 | 11.04 | 24.45 | 19.25 | 16.87 |
| Ground-Aware [47] | – | <span style="color:blue">21.65</span> | 13.25 | 9.91 | <span style="color:red">29.81</span> | 17.98 | 13.08 |
| MonoFlex [100] | – | 19.94 | 13.89 | <span style="color:red">12.07</span> | 28.23 | <span style="color:blue">19.75</span> | <span style="color:blue">16.89</span> |
| GUP Net [49] | – | 20.11 | <span style="color:blue">14.20</span> | 11.77 | – | – | – |
| **DEVIANT (Ours)** | – | <span style="color:red">21.88</span> | <span style="color:red">14.46</span> | <span style="color:blue">11.89</span> | <span style="color:blue">29.65</span> | <span style="color:red">20.44</span> | <span style="color:red">17.43</span> |

Table 4: **Results on KITTI Test cyclists and pedestrians** (Cyc/Ped) at $IoU_{3D} \geq 0.5$. Previous results are from the leader-board or papers. [Key: <span style="color:red">**Best**</span>, <span style="color:blue">**Second Best**</span>]

| Method | Extra | Cyc $AP_{3D|R_{40}}[\%](\blacktriangle)$ | | | Ped $AP_{3D|R_{40}}[\%](\blacktriangle)$ | | |
|---|---|---|---|---|---|---|---|
| | | Easy | Mod | Hard | Easy | Mod | Hard |
| DDMP-3D [79] | Depth | 4.18 | 2.50 | 2.32 | 4.93 | 3.55 | 3.01 |
| DFR-Net [105] | Depth | 5.69 | 3.58 | 3.10 | 6.09 | 3.62 | 3.39 |
| MonoDistill [14] | Depth | 5.53 | 2.81 | 2.40 | 12.79 | 8.17 | 7.45 |
| CaDDN [62] | LiDAR | 7.00 | 3.41 | 3.30 | 12.87 | 8.14 | 6.76 |
| DD3D [57] | LiDAR | 2.39 | 1.52 | 1.31 | 13.91 | 9.30 | 8.05 |
| MonoEF [103] | Odometry | 1.80 | 0.92 | 0.71 | 4.27 | 2.79 | 2.21 |
| MonoDIS-M [68] | – | 1.17 | 0.54 | 0.48 | 7.79 | 5.14 | 4.42 |
| MonoFlex [100] | – | 3.39 | 2.10 | 1.67 | 11.89 | 8.16 | 6.81 |
| GUP Net [49] | – | <span style="color:blue">4.18</span> | <span style="color:blue">2.65</span> | <span style="color:blue">2.09</span> | <span style="color:red">14.72</span> | <span style="color:red">9.53</span> | <span style="color:red">7.87</span> |
| **DEVIANT (Ours)** | – | <span style="color:red">5.05</span> | <span style="color:red">3.13</span> | <span style="color:red">2.59</span> | <span style="color:blue">13.43</span> | <span style="color:blue">8.65</span> | <span style="color:blue">7.69</span> |

MonoDistill [14], DEVIANT, an image-based method, also outperforms other methods that use extra information.
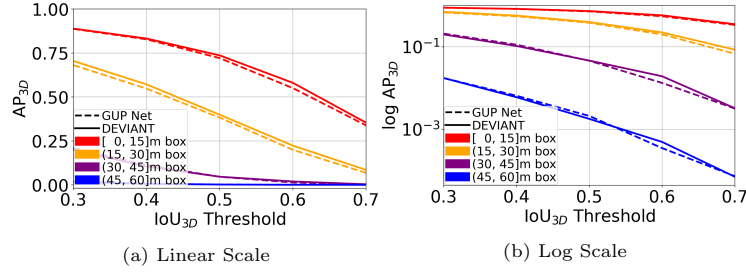
**Cyclists and Pedestrians.** Tab. 4 lists out the results of monocular 3D detection on KITTI Test Cyclist and Pedestrians. The results show that DEVIANT achieves SOTA results in the image-only category on the challenging Cyclists, and is competitive on Pedestrians.

### 5.2 KITTI Val Monocular 3D Detection

**Cars.** Tab. 5 summarizes the results of monocular 3D detection and BEV evaluation on KITTI Val split at two $IoU_{3D}$ thresholds of 0.7 and 0.5 [13,36]. We report the **median** model over 5 runs. The results show that DEVIANT outperforms the GUP Net [49] baseline by a significant margin. The biggest improvements shows up on the Easy set. Significant improvements are also on the Moderate and Hard sets. Interestingly, DEVIANT also outperforms DD3D [57] by a large margin when the large-dataset pretraining is not done (denoted by $DD3D^-$).

Table 5: **Results on KITTI Val cars**. Comparison with bigger CNN backbones in Tab. 16. [Key: **Best**, **Second Best**, ¯ = No pretrain]

| Method | Extra | $\mathrm{IoU_{3D}} \geq 0.7$ | | | | | | $\mathrm{IoU_{3D}} \geq 0.5$ | | | | | |
| | | $\mathrm{AP_{3D|R_{40}}}[\%](\uparrow)$ | | | $\mathrm{AP_{BEV|R_{40}}}[\%](\uparrow)$ | | | $\mathrm{AP_{3D|R_{40}}}[\%](\uparrow)$ | | | $\mathrm{AP_{BEV|R_{40}}}[\%](\uparrow)$ | | |
| | | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard |
| DDMP-3D [79] | Depth | 28.12 | 20.39 | 16.34 | – | – | – | – | – | – | – | – | – |
| PCT [80] | Depth | 38.39 | 27.53 | 24.44 | 47.16 | 34.65 | 28.47 | – | – | – | – | – | – |
| MonoDistill [14] | Depth | 24.31 | 18.47 | 15.76 | 33.09 | 25.40 | 22.16 | 65.69 | 49.35 | 43.49 | 71.45 | 53.11 | 46.94 |
| CaDDN [62] | LiDAR | 23.57 | 16.31 | 13.84 | – | – | – | – | – | – | – | – | – |
| PatchNet-C [69] | LiDAR | 24.51 | 17.03 | 13.25 | – | – | – | – | – | – | – | – | – |
| DD3D (DLA34) [57] | LiDAR | – | – | – | 33.5 | 26.0 | 22.6 | – | – | – | – | – | – |
| DD3D¯(DLA34) [57] | LiDAR | – | – | – | 26.8 | 20.2 | 16.7 | – | – | – | – | – | – |
| MonoEF [103] | Odometry | 18.26 | 16.30 | 15.24 | 26.07 | 25.21 | 21.61 | 57.98 | 51.80 | 49.34 | 63.40 | 61.13 | 53.22 |
| Kinematic [5] | Video | 19.76 | 14.10 | 10.47 | 27.83 | 19.72 | 15.10 | 55.44 | 39.47 | 31.26 | 61.79 | 44.68 | 34.56 |
| MonoRCNN [67] | – | 16.61 | 13.19 | 10.65 | 25.29 | 19.22 | 15.30 | – | – | – | – | – | – |
| MonoDLE [53] | – | 17.45 | 13.66 | 11.68 | 24.97 | 19.33 | 17.01 | 55.41 | 43.42 | 37.81 | 60.73 | 46.87 | 41.89 |
| GrooMeD-NMS [36] | – | 19.67 | 14.32 | 11.27 | 27.38 | 19.75 | 15.92 | 55.62 | 41.07 | 32.89 | 61.83 | 44.98 | 36.29 |
| Ground-Aware [47] | – | 23.63 | 16.16 | 12.06 | – | – | – | **60.92** | 42.18 | 32.02 | – | – | – |
| MonoFlex [100] | – | **23.64** | **17.51** | **14.83** | – | – | – | – | – | – | – | – | – |
| GUP Net (Reported) [49] | – | 22.76 | 16.46 | 13.72 | **31.07** | **22.94** | **19.75** | 57.62 | 42.33 | 37.59 | 61.78 | 47.06 | 40.88 |
| GUP Net (Retrained) [49] | – | 21.10 | 15.48 | 12.88 | 28.58 | 20.92 | 17.83 | 58.95 | **43.99** | **38.07** | **64.60** | **47.76** | **42.97** |
| **DEVIANT (Ours)** | – | **24.63** | **16.54** | **14.52** | **32.60** | **23.04** | **19.99** | **61.00** | **46.00** | **40.18** | **65.28** | **49.63** | **43.50** |



(a) Linear Scale      (b) Log Scale

Fig. 3: **AP₃D at different depths and IoU₃D thresholds** on KITTI Val Split.

Table 6: **Cross-dataset evaluation** of the KITTI Val model on KITTI Val and nuScenes frontal Val cars with depth MAE (↓). [Key: **Best**, **Second Best**]

| Method | KITTI Val | | | | nuScenes frontal Val | | | |
| | $0-20$ | $20-40$ | $40-\infty$ | All | $0-20$ | $20-40$ | $40-\infty$ | All |
| M3D-RPN [4] | 0.56 | 1.33 | 2.73 | 1.26 | 0.94 | 3.06 | 10.36 | 2.67 |
| MonoRCNN [67] | 0.46 | 1.27 | 2.59 | 1.14 | 0.94 | 2.84 | 8.65 | 2.39 |
| GUP Net [49] | **0.45** | **1.10** | **1.85** | **0.89** | **0.82** | **1.70** | **6.20** | **1.45** |
| **DEVIANT** | **0.40** | **1.09** | **1.80** | **0.87** | **0.76** | **1.60** | **4.50** | **1.26** |

**AP₃D at different depths and IoU₃D thresholds.** We next compare the AP₃D of DEVIANT and GUP Net in Fig. 3 at different distances in meters and IoU₃D matching criteria of $0.3 \rightarrow 0.7$ as in [36]. Fig. 3 shows that DEVIANT is effective over GUP Net [49] at all depths and higher IoU₃D thresholds.

**Cross-Dataset Evaluation.** Tab. 6 shows the result of our KITTI Val model on the KITTI Val and nuScenes [9] frontal Val images, using mean absolute error (MAE) of the depth of the boxes [67]. More details are in Appendix A3.1. DEVIANT outperforms GUP Net on most of the metrics on both the datasets, which confirms that DEVIANT generalizes better than CNNs. DEVIANT per-

Table 7: **Scale Augmentation vs Scale Equivariance** on KITTI Val cars. [Key: **Best**, Eqv= Equivariance, Aug= Augmentation]

| Method | Scale Eqv | Scale Aug | $AP_{3D|R_{40}}[\%](\blacktriangle)$ | | | $AP_{BEV|R_{40}}[\%](\blacktriangle)$ | | | $AP_{3D|R_{40}}[\%](\blacktriangle)$ | | | $AP_{BEV|R_{40}}[\%](\blacktriangle)$ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard |
| GUP Net [49] | | | 20.82 | 14.15 | 12.44 | 29.93 | 20.90 | 17.87 | **62.37** | 44.40 | 39.61 | **66.81** | 48.09 | 43.14 |
| | | ✓ | 21.10 | 15.48 | 12.88 | 28.58 | 20.92 | 17.83 | 58.95 | 43.99 | 38.07 | 64.60 | 47.76 | 42.97 |
| **DEVIANT** | ✓ | | 21.33 | 14.77 | 12.57 | 28.79 | 20.28 | 17.59 | 59.31 | 43.25 | 37.64 | 63.94 | 47.02 | 41.12 |
| | ✓ | ✓ | **24.63** | **16.54** | **14.52** | **32.60** | **23.04** | **19.99** | 61.00 | **46.00** | **40.18** | 65.28 | **49.63** | **43.50** |

Table 8: **Comparison of Equivariant Architectures** on KITTI Val cars. [Key: **Best**, Eqv= Equivariance, $^{\dagger}$= Retrained]

| Method | Eqv | $AP_{3D|R_{40}}[\%](\blacktriangle)$ | | | $AP_{BEV|R_{40}}[\%](\blacktriangle)$ | | | $AP_{3D|R_{40}}[\%](\blacktriangle)$ | | | $AP_{BEV|R_{40}}[\%](\blacktriangle)$ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard |
| DETR3D$^{\dagger}$ [85] | Learned | 1.94 | 1.26 | 1.09 | 4.41 | 3.06 | 2.79 | 20.09 | 13.80 | 12.78 | 26.51 | 18.49 | 17.36 |
| GUP Net [49] | 2D | 21.10 | 15.48 | 12.88 | 28.58 | 20.92 | 17.83 | 58.95 | 43.99 | 38.07 | 64.60 | 47.76 | 42.97 |
| **DEVIANT** | 2D+Depth | **24.63** | **16.54** | **14.52** | **32.60** | **23.04** | **19.99** | **61.00** | **46.00** | **40.18** | **65.28** | **49.63** | **43.50** |

Table 9: **Comparison with Dilated Convolution** on KITTI Val cars. [Key: **Best**]

| Method | Extra | $AP_{3D|R_{40}}[\%](\blacktriangle)$ | | | $AP_{BEV|R_{40}}[\%](\blacktriangle)$ | | | $AP_{3D|R_{40}}[\%](\blacktriangle)$ | | | $AP_{BEV|R_{40}}[\%](\blacktriangle)$ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard |
| D4LCN [20] | Depth | 22.32 | 16.20 | 12.30 | 31.53 | 22.58 | 17.87 | – | – | – | – | – | – |
| DCNN [97] | – | 21.66 | 15.49 | 12.90 | 30.22 | 22.06 | 19.01 | 57.54 | 43.12 | 38.80 | 63.29 | 46.86 | 42.42 |
| **DEVIANT** | – | **24.63** | **16.54** | **14.52** | **32.60** | **23.04** | **19.99** | **61.00** | **46.00** | **40.18** | **65.28** | **49.63** | **43.50** |

forms exceedingly well in the cross-dataset evaluation than [4, 49, 67]. We believe this happens because [4, 49, 67] rely on data or geometry to get the depth, while DEVIANT is equivariant to the depth translations, and therefore, outputs consistent depth. So, DEVIANT is more robust to data distribution changes.

**Alternatives to Equivariance.** We now compare with alternatives to equivariance in the following paragraphs.

**(a) Scale Augmentation.** A withstanding question in machine learning is the choice between equivariance and data augmentation [25]. Tab. 7 compares scale equivariance and scale augmentation. GUP Net [49] uses scale-augmentation and therefore, Tab. 7 shows that equivariance also benefits models which use scale-augmentation. This agrees with Tab. 2 of [74], where they observe that both augmentation and equivariance benefits classification on MNIST-scale dataset.

**(b) Other Equivariant Architectures.** We now benchmark adding depth (scale) equivariance to a 2D translation equivariant CNN and a transformer which learns the equivariance. Therefore, we compare DEVIANT with GUP Net [49] (a CNN), and DETR3D [85] (a transformer) in Tab. 8. As DETR3D does not report KITTI results, we trained DETR3D on KITTI using their public code. DEVIANT outperforms GUP Net and also surpasses DETR3D by a large margin. This happens because learning equivariance requires more data [90] compared to architectures which hardcode equivariance like CNN or DEVIANT.

**(c) Dilated Convolution.** DEVIANT adjusts the receptive field based on the object scale, and so, we compare with the dilated CNN (DCNN) [97] and D4LCN
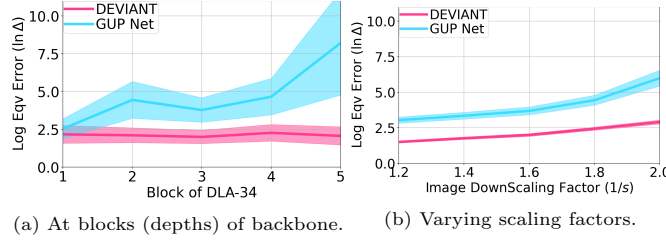
(a) At blocks (depths) of backbone.     (b) Varying scaling factors.

Fig. 4: **Log Equivariance Error** ($\Delta$) comparison for DEVIANT and GUP Net at **(a)** different blocks with random image scaling factors **(b)** different image scaling factors at depth 3. DEVIANT shows **lower** scale equivariance error than vanilla GUP Net [49].

[20] in Tab. 9. The results show that DCNN performs sub-par to DEVIANT. This is expected because dilation corresponds to integer scales [92] while the scaling is generally a float in monocular detection. D4LCN [20] uses monocular depth as input to adjust the receptive field. DEVIANT (without depth) also outperforms D4LCN on Hard cars, which are more distant.

**(d) Other Convolutions.** We now compare with other known convolutions in literature such as Log-polar convolution [106], Dilated convolution [97] convolution and DISCO [72] in Tab. 11. The results show that the log-polar convolution does not work well, and SES convolutions are better suited to embed depth (scale) equivariance. As described in Sec. 4, we investigate the behavior of log-polar convolution through a small experiment. We calculate the SSIM [86] of the original image and the image obtained after the upscaling, log-polar, inverse log-polar, and downscaling blocks. We then average the SSIM over all KITTI Val images. We repeat this experiment for multiple image heights and scaling factors. The ideal SSIM should have been one. However, Fig. 2c shows that SSIM does not reach 1 even after upscaling by 4. This result confirms that log-polar convolution loses information at low resolutions resulting in inaccurate detection.

Next, the results show that dilated convolution [97] performs sub-par to DEVIANT. Moreover, DISCO [72] also does not outperform SES convolution which agrees with the 2D tracking results of [72].

**(e) Feature Pyramid Network (FPN).** Our baseline GUP Net [49] uses FPN [44] and Tab. 5 shows that DEVIANT outperforms GUP Net. Hence, we conclude that equivariance also benefits models which use FPN.

**Comparison of Equivariance Error.** We next quantitatively evaluate the scale equivariance of DEVIANT vs. GUP Net [49], using the equivariance error metric [74]. The equivariance error $\Delta$ is the normalized difference between the scaled feature map and the feature map of the scaled image, and is given by $\Delta = \frac{1}{N} \sum_{i=1}^{N} \frac{||\mathcal{T}_{s_i} \Phi(h_i) - \Phi(\mathcal{T}_{s_i} h_i)||_2^2}{||\mathcal{T}_{s_i} \Phi(h_i)||_2^2}$, where $\Phi$ denotes the neural network, $\mathcal{T}_{s_i}$ is the scaling transformation for the image $i$, and $N$ is the total number of images. The equivariance error is zero if the scale equivariance is perfect. We plot the log of this error at different blocks of DEVIANT and GUP Net backbones and also plot at different downscaling of KITTI Val images in Fig. 4. The plots show that DEVIANT has low equivariance error than GUP Net. This is expected since the

Table 10: **3D and 2D detection** on KITTI Val cars.

| Method | IoU $\geq 0.7$ | | | | | | IoU $\geq 0.5$ | | | | | |
| | $\text{AP}_{3\text{D}|R_{40}}[\%](\blacktriangle)$ | | | $\text{AP}_{2\text{D}|R_{40}}[\%](\blacktriangle)$ | | | $\text{AP}_{3\text{D}|R_{40}}[\%](\blacktriangle)$ | | | $\text{AP}_{2\text{D}|R_{40}}[\%](\blacktriangle)$ | | |
| | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard |
| GUP Net [49] | 21.10 | 15.48 | 12.88 | 96.78 | 88.87 | 79.02 | 58.95 | 43.99 | 38.07 | 99.52 | 91.89 | 81.99 |
| **DEVIANT (Ours)** | 24.63 | 16.54 | 14.52 | 96.68 | 88.66 | 78.87 | 61.00 | 46.00 | 40.18 | 97.12 | 91.77 | 81.93 |

Table 11: **Ablation studies** on KITTI Val cars.

| Change from DEVIANT : | | $\text{IoU}_{3\text{D}} \geq 0.7$ | | | | | | $\text{IoU}_{3\text{D}} \geq 0.5$ | | | | | |
| Changed | From $\longrightarrow$ To | $\text{AP}_{3\text{D}|R_{40}}[\%](\blacktriangle)$ | | | $\text{AP}_{\text{BEV}|R_{40}}[\%](\blacktriangle)$ | | | $\text{AP}_{3\text{D}|R_{40}}[\%](\blacktriangle)$ | | | $\text{AP}_{\text{BEV}|R_{40}}[\%](\blacktriangle)$ | | |
| | | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard |
| | SES$\rightarrow$Vanilla | 21.10 | 15.48 | 12.88 | 28.58 | 20.92 | 17.83 | 58.95 | 43.99 | 38.07 | 64.60 | 47.76 | 42.97 |
| Convolution | SES$\rightarrow$Log-polar [106] | 9.19 | 6.77 | 5.78 | 16.39 | 11.15 | 9.80 | 40.51 | 27.62 | 23.90 | 45.66 | 31.34 | 25.80 |
| | SES$\rightarrow$Dilated [97] | 21.66 | 15.49 | 12.90 | 30.22 | 22.06 | 19.01 | 57.54 | 43.12 | 38.80 | 63.29 | 46.86 | 42.42 |
| | SES$\rightarrow$DISCO [72] | 20.21 | 13.84 | 11.46 | 28.56 | 19.38 | 16.41 | 55.22 | 39.76 | 35.37 | 59.46 | 43.16 | 38.52 |
| Downscale | 10% $\rightarrow$ 5% | 24.24 | 16.51 | 14.43 | 31.94 | 22.86 | 19.82 | 60.64 | 44.46 | 40.02 | 64.68 | 49.30 | 43.49 |
| $\alpha$ | 10% $\rightarrow$ 20% | 22.19 | 15.85 | 13.48 | 31.15 | 23.01 | 19.90 | 61.24 | 44.93 | 40.22 | 67.46 | 50.10 | 43.83 |
| BNP | SE$\rightarrow$ Vanilla | 24.39 | 16.20 | 14.36 | 32.43 | 22.53 | 19.70 | 62.81 | 46.14 | 40.38 | 67.87 | 50.23 | 44.08 |
| Scales | 3 $\rightarrow$ 1 | 23.20 | 16.29 | 13.63 | 31.76 | 23.23 | 19.97 | 61.90 | 46.66 | 40.61 | 67.37 | 50.31 | 43.93 |
| | 3 $\rightarrow$ 2 | 24.15 | 16.48 | 14.55 | 32.42 | 23.17 | 20.07 | 61.05 | 46.34 | 40.46 | 67.36 | 50.32 | 44.07 |
| — | **DEVIANT (best)** | **24.63** | **16.54** | **14.52** | **32.60** | **23.04** | **19.99** | **61.00** | **46.00** | **40.18** | **65.28** | **49.63** | **43.50** |

feature maps of the proposed DEVIANT are additionally equivariant to scale transformations (depth translations). We also visualize the equivariance error for a validation image and for the objects of this image in Fig. 12 in the supplementary. The qualitative plots also show a lower error for the proposed DEVIANT, which agrees with Fig. 4. Fig. 12a shows that equivariance error is particularly low for nearby cars which also justifies the good performance of DEVIANT on Easy (nearby) cars in Tabs. 3 and 5.

**Does 2D Detection Suffer?** We now investigate whether 2D detection suffers from using DEVIANT backbones in Tab. 10. The results show that DEVIANT introduces minimal decrease in the 2D detection performance. This is consistent with [73], who report that 2D tracking improves with the SE networks.

**Ablation Studies.** Tab. 11 compares the modifications of our approach on KITTI Val cars based on the experimental settings of Sec. 5.

**(a) Floating or Integer Downscaling?** We next investigate the question that whether one should use floating or integer downscaling factors for DEVIANT. We vary the downscaling factors as $(1{+}2\alpha, 1{+}\alpha, 1)$ and therefore, our scaling factor $s = \left(\frac{1}{1+2\alpha}, \frac{1}{1+\alpha}, 1\right)$. We find that $\alpha$ of 10% works the best. We again bring up the dilated convolution (Dilated) results at this point because dilation is a scale equivariant operation for integer downscaling factors [92] ($\alpha = 100\%, s = 0.5$). Tab. 11 results suggest that the downscaling factors should be floating numbers.

**(b) SE BNP.** As described in Sec. 4, we ablate DEVIANT against the case when only convolutions are SE but BNP layers are not. So, we place Scale-Projection [74] immediately after every SES convolution. Tab. 11 shows that such a network performs slightly sub-optimal to our final model.

**(c) Number of Scales.** We next ablate against the usage of Hermite scales. Using three scales performs better than using only one scale especially on Mod and Hard objects, and slightly better than using two scales.

Table 12: **Results on Waymo Val vehicles**. [Key: <span style="color:red">**Best**</span>, <span style="color:blue">**Second Best**</span>]

| IoU$_{3D}$ | Difficulty | Method | Extra | AP$_{3D}$ [%]($\uparrow$) | | | | APH$_{3D}$ [%]($\uparrow$) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | All | 0-30 | 30-50 | 50-∞ | All | 0-30 | 30-50 | 50-∞ |
| 0.7 | Level_1 | CaDDN [62] | LiDAR | 5.03 | 14.54 | 1.47 | 0.10 | 4.99 | 14.43 | 1.45 | 0.10 |
| | | PatchNet [50] in [80] | Depth | 0.39 | 1.67 | 0.13 | 0.03 | 0.39 | 1.63 | 0.12 | 0.03 |
| | | PCT [80] | Depth | 0.89 | 3.18 | 0.27 | 0.07 | 0.88 | 3.15 | 0.27 | 0.07 |
| | | M3D-RPN [4] in [62] | − | 0.35 | 1.12 | 0.18 | 0.02 | 0.34 | 1.10 | 0.18 | 0.02 |
| | | GUP Net (Retrained) [49] | − | <span style="color:blue">2.28</span> | <span style="color:blue">6.15</span> | <span style="color:blue">0.81</span> | <span style="color:red">0.03</span> | <span style="color:blue">2.27</span> | <span style="color:blue">6.11</span> | <span style="color:blue">0.80</span> | <span style="color:red">0.03</span> |
| | | **DEVIANT (Ours)** | − | <span style="color:red">2.69</span> | <span style="color:red">6.95</span> | <span style="color:red">0.99</span> | 0.02 | <span style="color:red">2.67</span> | <span style="color:red">6.90</span> | <span style="color:red">0.98</span> | 0.02 |
| 0.7 | Level_2 | CaDDN [62] | LiDAR | 4.49 | 14.50 | 1.42 | 0.09 | 4.45 | 14.38 | 1.41 | 0.09 |
| | | PatchNet [50] in [80] | Depth | 0.38 | 1.67 | 0.13 | 0.03 | 0.36 | 1.63 | 0.11 | 0.03 |
| | | PCT [80] | Depth | 0.66 | 3.18 | 0.27 | 0.07 | 0.66 | 3.15 | 0.26 | 0.07 |
| | | M3D-RPN [4] in [62] | − | 0.33 | 1.12 | 0.18 | <span style="color:red">0.02</span> | 0.33 | 1.10 | 0.17 | <span style="color:red">0.02</span> |
| | | GUP Net (Retrained) [49] | − | <span style="color:blue">2.14</span> | <span style="color:blue">6.13</span> | <span style="color:blue">0.78</span> | 0.02 | <span style="color:blue">2.12</span> | <span style="color:blue">6.08</span> | <span style="color:blue">0.77</span> | 0.02 |
| | | **DEVIANT (Ours)** | − | <span style="color:red">2.52</span> | <span style="color:red">6.93</span> | <span style="color:red">0.95</span> | 0.02 | <span style="color:red">2.50</span> | <span style="color:red">6.87</span> | <span style="color:red">0.94</span> | 0.02 |
| 0.5 | Level_1 | CaDDN [62] | LiDAR | 17.54 | 45.00 | 9.24 | 0.64 | 17.31 | 44.46 | 9.11 | 0.62 |
| | | PatchNet [50] in [80] | Depth | 2.92 | 10.03 | 1.09 | 0.23 | 2.74 | 9.75 | 0.96 | 0.18 |
| | | PCT [80] | Depth | 4.20 | 14.70 | 1.78 | 0.39 | 4.15 | 14.54 | 1.75 | 0.39 |
| | | M3D-RPN [4] in [62] | − | 3.79 | 11.14 | 2.16 | <span style="color:red">0.26</span> | 3.63 | 10.70 | 2.09 | <span style="color:blue">0.21</span> |
| | | GUP Net (Retrained) [49] | − | <span style="color:blue">10.02</span> | <span style="color:blue">24.78</span> | <span style="color:blue">4.84</span> | <span style="color:blue">0.22</span> | <span style="color:blue">9.94</span> | <span style="color:blue">24.59</span> | <span style="color:blue">4.78</span> | <span style="color:red">0.22</span> |
| | | **DEVIANT (Ours)** | − | <span style="color:red">10.98</span> | <span style="color:red">26.85</span> | <span style="color:red">5.13</span> | 0.18 | <span style="color:red">10.89</span> | <span style="color:red">26.64</span> | <span style="color:red">5.08</span> | 0.18 |
| 0.5 | Level_2 | CaDDN [62] | LiDAR | 16.51 | 44.87 | 8.99 | 0.58 | 16.28 | 44.33 | 8.86 | 0.55 |
| | | PatchNet [50] in [80] | Depth | 2.42 | 10.01 | 1.07 | 0.22 | 2.28 | 9.73 | 0.97 | 0.16 |
| | | PCT [80] | Depth | 4.03 | 14.67 | 1.74 | 0.36 | 4.15 | 14.51 | 1.71 | 0.35 |
| | | M3D-RPN [4] in [62] | − | 3.61 | 11.12 | 2.12 | <span style="color:red">0.24</span> | 3.46 | 10.67 | 2.04 | <span style="color:red">0.20</span> |
| | | GUP Net (Retrained) [49] | − | <span style="color:blue">9.39</span> | <span style="color:blue">24.69</span> | <span style="color:blue">4.67</span> | <span style="color:blue">0.19</span> | <span style="color:blue">9.31</span> | <span style="color:blue">24.50</span> | <span style="color:blue">4.62</span> | <span style="color:blue">0.19</span> |
| | | **DEVIANT (Ours)** | − | <span style="color:red">10.29</span> | <span style="color:red">26.75</span> | <span style="color:red">4.95</span> | 0.16 | <span style="color:red">10.20</span> | <span style="color:red">26.54</span> | <span style="color:red">4.90</span> | 0.16 |

### 5.3   Waymo Val Monocular 3D Detection

We also benchmark our method on the Waymo dataset [75] which has more variability than KITTI. Tab. 12 shows the results on Waymo Val split. The results show that DEVIANT outperforms the baseline GUP Net [49] on multiple levels and multiple thresholds. The biggest gains are on the nearby objects which is consistent with Tabs. 3 and 5. Interestingly, DEVIANT also outperforms PatchNet [50] and PCT [80] without using depth. Although the performance of DEVIANT lags CaDDN [62], it is important to stress that CaDDN uses LiDAR data in training, while DEVIANT is an image-only method.

## 6   Conclusions

This paper studies the modeling error in monocular 3D detection in detail and takes the first step towards convolutions equivariant to arbitrary 3D translations in the projective manifold. Since the depth is the hardest to estimate for this task, this paper proposes Depth EquiVarIAnt NeTwork (DEVIANT) built with existing scale equivariant steerable blocks. As a result, DEVIANT is equivariant to the depth translations in the projective manifold whereas vanilla networks are not. The additional depth equivariance forces the DEVIANT to learn consistent depth estimates and therefore, DEVIANT achieves SOTA detection results on KITTI and Waymo datasets in the image-only category and performs competitively to methods using extra information. Moreover, DEVIANT works better than vanilla networks in cross-dataset evaluation. Future works include applying the idea to Pseudo-LiDAR [83], and monocular 3D tracking.

# References

1. The KITTI Vision Benchmark Suite. http://www.cvlibs.net/datasets/kitti/eval_object.php?obj_benchmark=3d, accessed: 2022-07-03 8
2. Alhaija, H., Mustikovela, S., Mescheder, L., Geiger, A., Rother, C.: Augmented reality meets computer vision: Efficient data generation for urban driving scenes. IJCV (2018) 1
3. Bochkovskiy, A., Wang, C.Y., Liao, H.Y.M.: YOLOv4: Optimal speed and accuracy of object detection. arXiv preprint arXiv:2004.10934 (2020) 2
4. Brazil, G., Liu, X.: M3D-RPN: Monocular 3D region proposal network for object detection. In: ICCV (2019) 1, 2, 4, 10, 11, 14, 27, 34, 35
5. Brazil, G., Pons-Moll, G., Liu, X., Schiele, B.: Kinematic 3D object detection in monocular video. In: ECCV (2020) 1, 4, 9, 10, 24, 27, 36
6. Bronstein, M.: Convolution from first principles. https://towardsdatascience.com/deriving-convolution-from-first-principles-4ff124888028, accessed: 2021-08-13 1, 3, 4, 21
7. Bronstein, M., Bruna, J., Cohen, T., Veličković, P.: Geometric deep learning: Grids, groups, graphs, geodesics, and gauges. arXiv preprint arXiv:2104.13478 (2021) 3, 4, 21
8. Burns, B., Weiss, R., Riseman, E.: The non-existence of general-case view-invariants. In: Geometric invariance in computer vision (1992) 4, 5, 21, 22
9. Caesar, H., Bankiti, V., Lang, A., Vora, S., Liong, V., Xu, Q., Krishnan, A., Pan, Y., Baldan, G., Beijbom, O.: nuScenes: A multimodal dataset for autonomous driving. In: CVPR (2020) 8, 10, 34
10. Chabot, F., Chaouch, M., Rabarisoa, J., Teuliere, C., Chateau, T.: Deep MANTA: A coarse-to-fine many-task network for joint 2D and 3D vehicle analysis from monocular image. In: CVPR (2017) 4
11. Chen, X., Kundu, K., Zhang, Z., Ma, H., Fidler, S., Urtasun, R.: Monocular 3D object detection for autonomous driving. In: CVPR (2016) 4
12. Chen, X., Kundu, K., Zhu, Y., Berneshawi, A., Ma, H., Fidler, S., Urtasun, R.: 3D object proposals for accurate object class detection. In: NeurIPS (2015) 8
13. Chen, Y., Tai, L., Sun, K., Li, M.: MonoPair: Monocular 3D object detection using pairwise spatial relationships. In: CVPR (2020) 1, 4, 9
14. Chong, Z., Ma, X., Zhang, H., Yue, Y., Li, H., Wang, Z., Ouyang, W.: MonoDistill: Learning spatial features for monocular 3D object detection. In: ICLR (2022) 9, 10, 37
15. Cohen, T., Geiger, M., Köhler, J., Welling, M.: Spherical CNNs. In: ICLR (2018) 2, 3
16. Cohen, T., Welling, M.: Learning the irreducible representations of commutative lie groups. In: ICML (2014) 3
17. Cohen, T., Welling, M.: Group equivariant convolutional networks. In: ICML (2016) 3, 21
18. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: ImageNet: A large-scale hierarchical image database. In: CVPR (2009) 30
19. Dieleman, S., De Fauw, J., Kavukcuoglu, K.: Exploiting cyclic symmetry in convolutional neural networks. In: ICML (2016) 3, 21
20. Ding, M., Huo, Y., Yi, H., Wang, Z., Shi, J., Lu, Z., Luo, P.: Learning depth-guided convolutions for monocular 3D object detection. In: CVPR Workshops (2020) 4, 11, 12, 27

21. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. In: ICLR (2021) 3
22. Esteves, C., Allen-Blanchette, C., Zhou, X., Daniilidis, K.: Polar transformer networks. In: ICLR (2018) 3, 4
23. Fidler, S., Dickinson, S., Urtasun, R.: 3D object detection and viewpoint estimation with a deformable 3D cuboid model. In: NeurIPS (2012) 4
24. Freeman, W., Adelson, E.: The design and use of steerable filters. TPAMI (1991) 4, 7
25. Gandikota, K., Geiping, J., Lähner, Z., Czapliński, A., Moeller, M.: Training or architecture? how to incorporate invariance in neural networks. arXiv preprint arXiv:2106.10044 (2021) 2, 11, 21
26. Ganea, O.E., Bécigneul, G., Hofmann, T.: Hyperbolic neural networks. In: NeurIPS (2017) 2, 3
27. Geiger, A., Lenz, P., Stiller, C., Urtasun, R.: Vision meets robotics: The KITTI dataset. IJRR (2013) 38
28. Geiger, A., Lenz, P., Urtasun, R.: Are we ready for autonomous driving? the KITTI vision benchmark suite. In: CVPR (2012) 8
29. Ghosh, R., Gupta, A.: Scale steerable filters for locally scale-invariant convolutional neural networks. In: ICML Workshops (2019) 3, 4, 6, 7, 29
30. Hartley, R., Zisserman, A.: Multiple view geometry in computer vision. Cambridge university press (2003) 3, 4, 5, 23, 39
31. Henriques, J., Vedaldi, A.: Warped convolutions: Efficient invariance to spatial transformations. In: ICML (2017) 3
32. Jansson, Y., Lindeberg, T.: Scale-invariant scale-channel networks: Deep networks that generalise to previously unseen scales. IJCV (2021) 3, 4, 6
33. Jing, L.: Physical symmetry enhanced neural networks. Ph.D. thesis, Massachusetts Institute of Technology (2020) 3
34. Kanazawa, A., Sharma, A., Jacobs, D.: Locally scale-invariant convolutional neural networks. In: NeurIPS Workshops (2014) 3
35. Kingma, D., Ba, J.: Adam: A method for stochastic optimization. In: ICLR (2015) 31
36. Kumar, A., Brazil, G., Liu, X.: GrooMeD-NMS: Grouped mathematically differentiable NMS for monocular 3D object detection. In: CVPR (2021) 4, 7, 9, 10, 32, 34, 37, 38
37. Kumar, A., Marks, T., Mou, W., Wang, Y., Jones, M., Cherian, A., Koike-Akino, T., Liu, X., Feng, C.: LUVLi face alignment: Estimating landmarks' location, uncertainty, and visibility likelihood. In: CVPR (2020) 4
38. Kumar, A., Prabhakaran, V.: Estimation of bandlimited signals from the signs of noisy samples. In: ICASSP (2013) 7
39. Lambert, J., Liu, Z., Sener, O., Hays, J., Koltun, V.: MSeg: A composite dataset for multi-domain semantic segmentation. In: CVPR (2020) 36
40. LeCun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. Proceedings of the IEEE (1998) 3, 4
41. Lee, J., Han, M., Ko, D., Suh, I.: From big to small: Multi-scale local planar guidance for monocular depth estimation. arXiv preprint arXiv:1907.10326 (2019) 35, 36, 39
42. Li, P., Zhao, H., Liu, P., Cao, F.: RTM3D: Real-time monocular 3D detection from object keypoints for autonomous driving. In: ECCV (2020) 1, 4

43. Lian, Q., Ye, B., Xu, R., Yao, W., Zhang, T.: Geometry-aware data augmentation for monocular 3D object detection. arXiv preprint arXiv:2104.05858 (2021) 2, 4
44. Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: CVPR (2017) 12, 30
45. Liu, L., Lu, J., Xu, C., Tian, Q., Zhou, J.: Deep fitting degree scoring network for monocular 3D object detection. In: CVPR (2019) 1, 4
46. Liu, X., Xue, N., Wu, T.: Learning auxiliary monocular contexts helps monocular 3D object detection. In: AAAI (2022) 4
47. Liu, Y., Yixuan, Y., Liu, M.: Ground-aware monocular 3D object detection for autonomous driving. Robotics and Automation Letters (2021) 2, 9, 10
48. Liu, Z., Zhou, D., Lu, F., Fang, J., Zhang, L.: AutoShape: Real-time shape-aware monocular 3D object detection. In: ICCV (2021) 4, 9
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 (2021) 1, 2, 4, 7, 8, 9, 10, 11, 12, 13, 14, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42
50. Ma, X., Liu, S., Xia, Z., Zhang, H., Zeng, X., Ouyang, W.: Rethinking Pseudo-LiDAR representation. In: ECCV (2020) 4, 14
51. Ma, X., Ouyang, W., Simonelli, A., Ricci, E.: 3D object detection from images for autonomous driving: A survey. arXiv preprint arXiv:2202.02980 (2022) 4
52. 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: ICCV (2019) 4
53. 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 (2021) 2, 8, 10, 25, 27, 28, 32
54. Marcos, D., Kellenberger, B., Lobry, S., Tuia, D.: Scale equivariance in CNNs with vector fields. In: ICML Workshops (2018) 3
55. Marcos, D., Volpi, M., Komodakis, N., Tuia, D.: Rotation equivariant vector field networks. In: ICCV (2017) 3
56. Micheli, A.: Neural network for graphs: A contextual constructive approach. IEEE Transactions on Neural Networks (2009) 3
57. Park, D., Ambrus, R., Guizilini, V., Li, J., Gaidon, A.: Is Pseudo-LiDAR needed for monocular 3D object detection? In: ICCV (2021) 4, 8, 9, 10, 33, 37
58. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: PyTorch: An imperative style, high-performance deep learning library. In: NeurIPS (2019) 8
59. Payet, N., Todorovic, S.: From contours to 3D object detection and pose estimation. In: ICCV (2011) 4
60. Pepik, B., Stark, M., Gehler, P., Schiele, B.: Multi-view and 3D deformable part models. TPAMI (2015) 4
61. Rath, M., Condurache, A.: Boosting deep neural networks with geometrical prior knowledge: A survey. arXiv preprint arXiv:2006.16867 (2020) 1, 3, 4, 21
62. Reading, C., Harakeh, A., Chae, J., Waslander, S.: Categorical depth distribution network for monocular 3D object detection. In: CVPR (2021) 4, 8, 9, 10, 14, 30, 32, 37
63. Rematas, K., Kemelmacher-Shlizerman, I., Curless, B., Seitz, S.: Soccer on your tabletop. In: CVPR (2018) 1

64. Ren, S., He, K., Girshick, R., Sun, J.: Faster R-CNN: Towards real-time object detection with region proposal networks. In: NeurIPS (2015) 1
65. Saxena, A., Driemeyer, J., Ng, A.: Robotic grasping of novel objects using vision. IJRR (2008) 1
66. Shi, S., Wang, X., Li, H.: PointRCNN: 3D object proposal generation and detection from point cloud. In: CVPR (2019) 4
67. Shi, X., Ye, Q., Chen, X., Chen, C., Chen, Z., Kim, T.K.: Geometry-based distance decomposition for monocular 3D object detection. In: ICCV (2021) 2, 8, 9, 10, 11, 34, 35, 37, 39
68. Simonelli, A., Bulò, S., Porzi, L., Antequera, M., Kontschieder, P.: Disentangling monocular 3D object detection: From single to multi-class recognition. TPAMI (2020) 1, 8, 9, 34, 37
69. Simonelli, A., Bulò, S., Porzi, L., Kontschieder, P., Ricci, E.: Are we missing confidence in Pseudo-LiDAR methods for monocular 3D object detection? In: ICCV (2021) 4, 9, 10, 36
70. Simonelli, A., Bulò, S., Porzi, L., López-Antequera, M., Kontschieder, P.: Disentangling monocular 3D object detection. In: ICCV (2019) 8, 34
71. Simonelli, A., Bulò, S., Porzi, L., Ricci, E., Kontschieder, P.: Towards generalization across depth for monocular 3D object detection. In: ECCV (2020) 1, 2, 4
72. Sosnovik, I., Moskalev, A., Smeulders, A.: DISCO: accurate discrete scale convolutions. In: BMVC (2021) 12, 13
73. Sosnovik, I., Moskalev, A., Smeulders, A.: Scale equivariance improves siamese tracking. In: WACV (2021) 4, 6, 8, 13, 29, 30
74. Sosnovik, I., Szmaja, M., Smeulders, A.: Scale-equivariant steerable networks. In: ICLR (2020) 3, 4, 6, 7, 11, 12, 13, 26, 29, 37
75. Sun, P., Kretzschmar, H., Dotiwalla, X., Chouard, A., Patnaik, V., Tsui, P., Guo, J., Zhou, Y., Chai, Y., Caine, B., Vasudevan, V., Han, W., Ngiam, J., Zhao, H., Timofeev, A., Ettinger, S., Krivokon, M., Gao, A., Joshi, A., Zhang, Y., Shlens, J., Chen, Z., Anguelov, D.: Scalability in perception for autonomous driving: Waymo open dataset. In: CVPR (2020) 8, 14
76. Tang, Y., Dorn, S., Savani, C.: Center3D: Center-based monocular 3D object detection with joint depth understanding. arXiv preprint arXiv:2005.13423 (2020) 1, 2, 4
77. Thayalan-Vaz, S., M, S., Santhakumar, K., Ravi Kiran, B., Gauthier, T., Yogamani, S.: Exploring 2D data augmentation for 3D monocular object detection. arXiv preprint arXiv:2104.10786 (2021) 2, 4
78. Thomas, N., Smidt, T., Kearnes, S., Yang, L., Li, L., Kohlhoff, K., Riley, P.: Tensor field networks: Rotation-and translation-equivariant neural networks for 3D point clouds. arXiv preprint arXiv:1802.08219 (2018) 3
79. 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: CVPR (2021) 9, 10
80. Wang, L., Zhang, L., Zhu, Y., Zhang, Z., He, T., Li, M., Xue, X.: Progressive coordinate transforms for monocular 3D object detection. In: NeurIPS (2021) 8, 9, 10, 14, 37, 39
81. Wang, R., Walters, R., Yu, R.: Incorporating symmetry into deep dynamics models for improved generalization. In: ICLR (2021) 3
82. Wang, X., Zhang, S., Yu, Z., Feng, L., Zhang, W.: Scale-equalizing pyramid convolution for object detection. In: CVPR (2020) 33

83. Wang, Y., Chao, W.L., Garg, D., Hariharan, B., Campbell, M., Weinberger, K.: Pseudo-LiDAR from visual depth estimation: Bridging the gap in 3D object detection for autonomous driving. In: CVPR (2019) 4, 14

84. Wang, Y., Chen, X., You, Y., Li, L., Hariharan, B., Campbell, M., Weinberger, K., Chao, W.L.: Train in Germany, test in the USA: Making 3D object detectors generalize. In: CVPR (2020) 35

85. Wang, Y., Guizilini, V., Zhang, T., Wang, Y., Zhao, H., Solomon, J.: DETR3D: 3D object detection from multi-view images via 3D-to-2D queries. In: CoRL (2021) 11, 39

86. Wang, Z., Bovik, A., Sheikh, H., Simoncelli, E.: Image quality assessment: from error visibility to structural similarity. TIP (2004) 12

87. Weiler, M., Forré, P., Verlinde, E., Welling, M.: Coordinate independent convolutional networks–isometry and gauge equivariant convolutions on riemannian manifolds. arXiv preprint arXiv:2106.06020 (2021) 3

88. Weiler, M., Hamprecht, F., Storath, M.: Learning steerable filters for rotation equivariant CNNs. In: CVPR (2018) 3

89. Wilk, M.v.d., Bauer, M., John, S., Hensman, J.: Learning invariances using the marginal likelihood. In: NeurIPS (2018) 3

90. Worrall, D., Brostow, G.: Cubenet: Equivariance to 3D rotation and translation. In: ECCV (2018) 3, 4, 11

91. Worrall, D., Garbin, S., Turmukhambetov, D., Brostow, G.: Harmonic networks: Deep translation and rotation equivariance. In: CVPR (2017) 3

92. Worrall, D., Welling, M.: Deep scale-spaces: Equivariance over scale. In: NeurIPS (2019) 4, 12, 13, 27

93. Wu, Y., Johnson, J.: Rethinking "batch" in batchnorm. arXiv preprint arXiv:2105.07576 (2021) 33

94. Xu, Y., Xiao, T., Zhang, J., Yang, K., Zhang, Z.: Scale-invariant convolutional neural networks. arXiv preprint arXiv:1411.6369 (2014) 4

95. Yang, G., Ramanan, D.: Upgrading optical flow to 3D scene flow through optical expansion. In: CVPR (2020) 8

96. Yeh, R., Hu, Y.T., Schwing, A.: Chirality nets for human pose regression. NeurIPS (2019) 3

97. Yu, F., Koltun, V.: Multi-scale context aggregation by dilated convolutions. In: ICLR (2015) 11, 12, 13, 27

98. Yu, F., Wang, D., Shelhamer, E., Darrell, T.: Deep layer aggregation. In: CVPR (2018) 30

99. Zhang, Y., Ma, X., Yi, S., Hou, J., Wang, Z., Ouyang, W., Xu, D.: Learning geometry-guided depth via projective modeling for monocular 3D object detection. arXiv preprint arXiv:2107.13931 (2021) 2

100. Zhang, Y., Lu, J., Zhou, J.: Objects are different: Flexible monocular 3D object detection. In: CVPR (2021) 1, 4, 9, 10, 37

101. Zhou, A., Knowles, T., Finn, C.: Meta-learning symmetries by reparameterization. In: ICLR (2021) 3

102. Zhou, X., Wang, D., Krähenbühl, P.: Objects as points. arXiv preprint arXiv:1904.07850 (2019) 1

103. Zhou, Y., He, Y., Zhu, H., Wang, C., Li, H., Jiang, Q.: MonoEF: Extrinsic parameter free monocular 3D object detection. TPAMI (2021) 2, 4, 9, 10, 37

104. Zhu, W., Qiu, Q., Calderbank, R., Sapiro, G., Cheng, X.: Scale-equivariant neural networks with decomposed convolutional filters. arXiv preprint arXiv:1909.11193 (2019) 3, 4, 6, 7, 37

105. 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 (2021) 9
106. Zwicke, P., Kiss, I.: A new implementation of the mellin transform and its application to radar classification of ships. TPAMI (1983) 2, 3, 7, 12, 13

# DEVIANT: Depth EquiVarIAnt NeTwork for Monocular 3D Object Detection Supplementary Material

## A1    Supportive Explanations

We now add some explanations which we could not put in the main paper because of the space constraints.

### A1.1    Equivariance vs Augmentation

Equivariance adds suitable inductive bias to the backbone [17, 19] and is not learnt. Augmentation adds transformations to the input data during training or inference.

Equivariance and data augmentation have their own pros and cons. Equivariance models the physics better, is mathematically principled and is so more agnostic to data distribution shift compared to the data augmentation. A downside of equivariance compared to the augmentation is equivariance requires mathematical modelling, may not always exist [8], is not so intuitive and generally requires more flops for inference. On the other hand, data augmentation is simple, intuitive and fast, but is not mathematically principled. The choice between equivariance and data augmentation is a withstanding question in machine learning [25].

### A1.2    Why do 2D CNN detectors generalize?

We now try to understand why 2D CNN detectors generalize well. Consider an image $h(u, v)$ and $\Phi$ be the CNN. Let $\mathcal{T}_\mathbf{t}$ denote the translation in the $(u, v)$ space. The 2D translation equivariance [6, 7, 61] of the CNN means that

$$\Phi(\mathcal{T}_\mathbf{t} h(u, v)) = \mathcal{T}_\mathbf{t} \Phi(h(u, v))$$
$$\implies \Phi(h(u + t_u, v + t_v)) = \Phi(h(u, v)) + (t_u, t_v) \tag{5}$$

where $(t_u, t_v)$ is the translation in the $(u, v)$ space.

Assume the CNN predicts the object position in the image as $(u', v')$. Then, we write

$$\Phi(h(u, v)) = (\hat{u}, \hat{v}) \tag{6}$$

Now, we want the CNN to predict the output the position of the same object translated by $(t_u, t_v)$. The new image is thus $h(u + t_u, v + t_v)$. The CNN easily predicts the translated position of the object because all CNN is to do is to invoke its 2D translation equivariance of Eq. (5), and translate the previous prediction by the same amount. In other words,

$$\Phi(h(u + t_u, v + t_v)) = \Phi(h(u, v)) + (t_u, t_v)$$

Fig. 5: **Equivariance exists** for the patch plane when there is depth translation of the ego camera. Downscaling converts image $h$ to image $h'$.



Fig. 6: **Example of non-existence of equivariance** [8] when there is $180°$ rotation of the ego camera. No transformation can convert image $h$ to image $h'$.

$$= (\hat{u}, \hat{v}) + (t_u, t_v)$$
$$= (\hat{u} + t_u, \hat{v} + t_v)$$

Intuitively, equivariance is a disentaglement method. The 2D translation equivariance disentangles the 2D translations $(t_u, t_v)$ from the original image $h$ and therefore, the network generalizes to unseen 2D translations.

### A1.3   Existence and Non-existence of Equivariance

The result from [8] says that generic projective equivariance does not exist in particular with rotation transformations. We now show an example of when the equivariance exists and does not exist in the projective manifold in Figs. 5 and 6 respectively.

### A1.4   Why do not Monocular 3D CNN detectors generalize?

Monocular 3D CNN detectors do not generalize well because they are not equivariant to arbitrary 3D translations in the projective manifold. To show this, let $H(x, y, z)$ denote a 3D point cloud. The monocular detection network $\Phi$ operates on the projection $h(u, v)$ of this point cloud $H$ to output the position $(\hat{x}, \hat{y}, \hat{z})$ as

$$\Phi(\mathbf{K}H(x, y, z)) = (\hat{x}, \hat{y}, \hat{z})$$

$$\implies \Phi(h(u,v)) = (\hat{x}, \hat{y}, \hat{z}),$$

where $\mathbf{K}$ denotes the projection operator. We translate this point cloud by an arbitrary 3D translation of $(t_x, t_y, t_z)$ to obtain the new point cloud $H(x+t_x, y+t_y, z+t_z)$. Then, we again ask the monocular detector $\Phi$ to do prediction over the translated point cloud. However, we find that

$$\Phi(\mathbf{K}H(x+t_x, y+t_y, z+t_z)) \neq \Phi(h(u + \mathbf{K}(t_x,t_y,t_z), v + \mathbf{K}(t_x,t_y,t_z)))$$
$$= \Phi(h(u,v)) + \mathbf{K}(t_x,t_y,t_z)$$
$$\implies \Phi(\mathbf{K}H(x+t_x, y+t_y, z+t_z)) \neq \Phi(\mathbf{K}H(x,y,z)) + \mathbf{K}(t_x,t_y,t_z)$$

In other words, the projection operator $\mathbf{K}$ does not distribute over the point cloud $H$ and arbitrary 3D translation of $(t_x, t_y, t_z)$. Hence, if the network $\Phi$ is a vanilla CNN (existing monocular backbone), it can no longer invoke its 2D translation equivariance of Eq. (5) to get the new 3D coordinates $(\hat{x} + t_x, \hat{y} + t_y, \hat{z} + t_z)$.

Note that the LiDAR based 3D detectors with 3D convolutions do not suffer from this problem because they do not involve any projection operator $\mathbf{K}$. Thus, this problem exists only in monocular 3D detection. This makes monocular 3D detection different from 2D and LiDAR based 3D object detection.

### A1.5   Overview of Theorem 1

We now pictorially provide the overview of Theorem 1 (Example 13.2 from [30]), which links the planarity and projective transformations in the continuous world in Fig. 7.



Fig. 7: **Overview** of Theorem 1 (Example 13.2 from [30]), which links the planarity and projective transformations in the continuous world.

### A1.6     Approximation of Corollary 1

We now give the approximation under which Corollary 1 is valid. We assume that the ego camera does not undergo any rotation. Hence, we substitute $\mathbf{R} = \mathbf{I}$ in Eq. (1) to get

$$h(u - u_0, v - v_0) = h' \left( f \frac{\left(1 + t_x \frac{m}{p}\right)(u - u_0) + t_x \frac{n}{p}(v - v_0) + t_x \frac{o}{p} f}{t_z \frac{m}{p}(u - u_0) + t_z \frac{n}{p}(v - v_0) + \left(1 + t_z \frac{o}{p}\right) f}, \right.$$
$$\left. f \frac{t_y \frac{m}{p}(u - u_0) + \left(1 + t_y \frac{n}{p}\right)(v - v_0) + t_y \frac{o}{p} f}{t_z \frac{m}{p}(u - u_0) + t_z \frac{n}{p}(v - v_0) + \left(1 + t_z \frac{o}{p}\right) f} \right). \qquad (7)$$

Next, we use the assumption that the ego vehicle moves in the $z$-direction as in [5], *i.e.*, substitute $t_x = t_y = 0$ to get

$$h(u - u_0, v - v_0) = h' \left( \frac{u - u_0}{\frac{t_z}{f} \frac{m}{p}(u - u_0) + \frac{t_z}{f} \frac{n}{p}(v - v_0) + \left(1 + t_z \frac{o}{p}\right)}, \right.$$
$$\left. \frac{v - v_0}{\frac{t_z}{f} \frac{m}{p}(u - u_0) + \frac{t_z}{f} \frac{n}{p}(v - v_0) + \left(1 + t_z \frac{o}{p}\right)} \right). \qquad (8)$$

The patch plane is $mx + ny + oz + p = 0$. We consider the planes in the front of camera. Without loss of generality, consider $p < 0$ and $o > 0$.

We first write the denominator $D$ of RHS term in Eq. (8) as

$$D = \frac{t_z}{f} \frac{m}{p}(u - u_0) + \frac{t_z}{f} \frac{n}{p}(v - v_0) + \left(1 + t_z \frac{o}{p}\right)$$
$$= 1 + \frac{t_z}{p} \left( \frac{m}{f}(u - u_0) + \frac{n}{f}(v - v_0) + o \right)$$

Because we considered patch planes in front of the camera, $p < 0$. Also consider $t_z < 0$, which implies $t_z/p > 0$. Now, we bound the term in the parantheses of the above equation as

$$D \leq 1 + \frac{t_z}{p} \left| \frac{m}{f}(u - u_0) + \frac{n}{f}(v - v_0) + o \right|$$
$$\leq 1 + \frac{t_z}{p} \left( \left| \frac{m}{f}(u - u_0) \right| + \left| \frac{n}{f}(v - v_0) \right| + |o| \right) \quad \text{by Triangle inequality}$$
$$\leq 1 + \frac{t_z}{p} \left( \frac{|m|}{f} \frac{W}{2} + \frac{|n|}{f} \frac{H}{2} + o \right), (u - u_0) \leq \frac{W}{2}, (v - v_0) \leq \frac{H}{2}, |o| = o$$
$$\leq 1 + \frac{t_z}{p} \left( \frac{|m|}{f} \frac{W}{2} + \frac{|n|}{f} \frac{W}{2} + o \right), H \leq W$$
$$\leq 1 + \frac{t_z}{p} \left( \frac{(|m| + |n|)W}{2f} + o \right),$$

If the coefficients of the patch plane $m, n, o$, its width $W$ and focal length $f$ follow the relationship $\frac{(|m|+|n|)W}{2f} << o$, the patch plane is "approximately" parallel to the image plane. Then, a few quantities can be ignored in the denominator $D$ to get

$$D \approx 1 + t_z \frac{o}{p} \qquad (9)$$

Therefore, the RHS of Eq. (8) gets simplified and we obtain

$$\mathcal{T}_s : h(u - u_0, v - v_0) \approx h' \left( \frac{u - u_0}{1 + t_z \frac{o}{p}}, \frac{v - v_0}{1 + t_z \frac{o}{p}} \right) \qquad (10)$$

An immediate benefit of using the approximation is Eq. (2) does not depend on the distance of the patch plane from the camera. This is different from wide-angle camera assumption, where the ego camera is assumed to be far from the patch plane. Moreover, patch planes need not be perfectly aligned with the image plane for Eq. (2). Even small enough perturbed patch planes work. We next show the approximation in the Fig. 8 with $\theta$ denoting the deviation from the perfect parallel plane. The deviation $\theta$ is about 3 degrees for the KITTI dataset while it is 6 degrees for the Waymo dataset.



Fig. 8: **Approximation of Corollary 1**. Bold shows the patch plane parallel to the image plane. The dotted line shows the approximated patch plane.

*e.g.* The following are valid patch planes for KITTI images whose focal length $f = 707$ and width $W = 1242$.

$$-0.05x + 0.05y + z = 30$$
$$0.05x - 0.05y + z = 30 \qquad (11)$$

The following are valid patch planes for Waymo images whose focal length $f = 2059$ and width $W = 1920$.

$$-0.1x + 0.1y + z = 30$$
$$0.1x - 0.1y + z = 30 \qquad (12)$$

Although the assumption is slightly restrictive, we believe our method shows improvements on both KITTI and Waymo datasets because the car patches are approximately parallel to image planes and also because the depth remains the hardest parameter to estimate [53].

### A1.7    Scale Equivariance of SES Convolution for Images

[74] derive the scale equivariance of SES convolution for a 1D signal. We simply follow on their footsteps to get the scale equivariance of SES convolution for a 2D image $h(u, v)$ for the sake of completeness. Let the scaling of the image $h$ be $s$. Let $*$ denote the standard vanilla convolution and $\Psi$ denote the convolution filter. Then, the convolution of the downscaled image $\mathcal{T}_s(h)$ with the filter $\Psi$ is given by

$$
\begin{aligned}
[\mathcal{T}_s(h) * \Psi]\,(u, v) & \\
&= \int \int h\left(\frac{u'}{s}, \frac{v'}{s}\right) \Psi(u' - u, v' - v)du'dv' \\
&= s^2 \int \int h\left(\frac{u'}{s}, \frac{v'}{s}\right) \Psi\left(s\frac{u' - u}{s}, s\frac{v' - v}{s}\right) d\left(\frac{u'}{s}\right) d\left(\frac{v'}{s}\right) \\
&= s^2 \int \int h\left(\frac{u'}{s}, \frac{v'}{s}\right) \mathcal{T}_{s^{-1}}\left[\Psi\left(\frac{u' - u}{s}, \frac{v' - v}{s}\right)\right] d\left(\frac{u'}{s}\right) d\left(\frac{v'}{s}\right) \\
&= s^2 \int \int h\left(\frac{u'}{s}, \frac{v'}{s}\right) \mathcal{T}_{s^{-1}}\left[\Psi\left(\frac{u'}{s} - \frac{u}{s}, \frac{v'}{s} - \frac{v}{s}\right)\right] d\left(\frac{u'}{s}\right) d\left(\frac{v'}{s}\right) \\
&= s^2 \left[h * \mathcal{T}_{s^{-1}}(\Psi)\right]\left(\frac{u}{s}, \frac{v}{s}\right) \\
&= s^2 \mathcal{T}_s\left[h * \mathcal{T}_{s^{-1}}(\Psi)\right](u, v).
\end{aligned}
\tag{13}
$$

Next, [74] re-parametrize the SES filters by writing $\Psi_\sigma(u, v) = \frac{1}{\sigma^2}\Psi\left(\frac{u}{\sigma}, \frac{v}{\sigma}\right)$. Substituting in Eq. (13), we get

$$
[\mathcal{T}_s(h) * \Psi_\sigma]\,(u, v) = s^2 \mathcal{T}_s\left[h * \mathcal{T}_{s^{-1}}(\Psi_\sigma)\right](u, v)
\tag{14}
$$

Moreover, the re-parametrized filters are separable [74] by construction and so, one can write

$$
\Psi_\sigma(u, v) = \Psi_\sigma(u)\Psi_\sigma(v).
\tag{15}
$$

The re-parametrization and separability leads to the important property that

$$
\begin{aligned}
\mathcal{T}_{s^{-1}}\left(\Psi_\sigma(u, v)\right) &= \mathcal{T}_{s^{-1}}\left(\Psi_\sigma(u)\Psi_\sigma(v)\right) \\
&= \mathcal{T}_{s^{-1}}\left(\Psi_\sigma(u)\right)\mathcal{T}_{s^{-1}}\left(\Psi_\sigma(v)\right) \\
&= s^{-2}\Psi_{s^{-1}\sigma}(u)\Psi_{s^{-1}\sigma}(v) \\
&= s^{-2}\Psi_{s^{-1}\sigma}(u, v).
\end{aligned}
\tag{16}
$$

Substituting above in the RHS of Eq. (14), we get

$$
\begin{aligned}
[\mathcal{T}_s(h) * \Psi_\sigma]\,(u, v) &= s^2 \mathcal{T}_s\left[h * s^{-2}\Psi_{s^{-1}\sigma}\right](u, v) \\
\implies [\mathcal{T}_s(h) * \Psi_\sigma]\,(u, v) &= \mathcal{T}_s\left[h * \Psi_{s^{-1}\sigma}\right](u, v),
\end{aligned}
\tag{17}
$$

which is a cleaner form of Eq. (13). Eq. (17) says that convolving the downscaled image with a filter is same as the downscaling the result of convolving the image with the upscaled filter [74]. This additional constraint regularizes the scale (depth) predictions for the image, leading to better generalization.

Table 13: **Comparison of Methods** on the basis of inputs, convolution kernels, outputs and whether output are scale-constrained.

| Method | Input Frame | #Conv Kernel | Output | Output Constrained for Scales? |
|---|---|---|---|---|
| Vanilla CNN | 1 | 1 | 4D | ✗ |
| Depth-Aware [4] | 1 | > 1 | 4D | ✗ |
| Dilated CNN [97] | 1 | > 1 | 5D | Integer [92] |
| **DEVIANT** | 1 | > 1 | 5D | Float |
| Depth-guided [20] | 1 + Depth | 1 | 4D | Integer [92] |
| Kinematic3D [5] | > 1 | 1 | 5D | ✗ |

## A1.8   Why does DEVIANT generalize better compared to CNN backbone?

DEVIANT models the physics better compared to the CNN backbone. CNN generalizes better for 2D detection because of the 2D translation equivariance in the Euclidean manifold. However, monocular 3D detection does not belong to the Euclidean manifold but is a task of the projective manifold. Modeling translation equivariance in the correct manifold improves generalization. For monocular 3D detection, we take the first step towards the general 3D translation equivariance by embedding equivariance to depth translations. The 3D depth equivariance in DEVIANT uses Eq. (14) and thus imposes an additional constraint on the feature maps. This additional constraint results in consistent depth estimates from the current image and a virtual image (obtained by translating the ego camera), and therefore, better generalization than CNNs. On the other hand, CNNs, by design, do not constrain the depth estimates from the current image and a virtual image (obtained by translating the ego camera), and thus, their depth estimates are entirely data-driven.

## A1.9   Why not Fixed Scale Assumption?

We now answer the question of keeping the fixed scale assumption. If we assume fixed scale assumption, then vanilla convolutional layers have the right equivariance. However, we do not keep this assumption because the ego camera translates along the depth in driving scenes and also, because the depth is the hardest parameter to estimate [53] for monocular detection. So, zero depth translation or fixed scale assumption is always violated.

## A1.10   Comparisons with Other Methods

We now list out the differences between different convolutions and monocular detection methods in Tab. 13. Kinematic3D [5] does not constrain the output at feature map level, but at system level using Kalman Filters. The closest to our method is the Dilated CNN (DCNN) [97]. We show in Tab. 9 that DEVIANT outperforms Dilated CNN.

### A1.11   Why is Depth the hardest among all parameters?

Images are the 2D projections of the 3D scene, and therefore, the depth is lost during projection. Recovering this depth is the most difficult to estimate, as shown in Tab. 1 of [53]. Monocular detection task involves estimating 3D center, 3D dimensions and the yaw angle. The right half of Tab. 1 in [53] shows that if the ground truth 3D center is replaced with the predicted center, the detection reaches a minimum. Hence, 3D center is the most difficult to estimate among center, dimensions and pose. Most monocular 3D detectors further decompose the 3D center into projected (2D) center and depth. Out of projected center and depth, Tab. 1 of [53] shows that replacing ground truth depth with the predicted depth leads to inferior detection compared to replacing ground truth projected center with the predicted projected center. Hence, we conclude that depth is the hardest parameter to estimate.

Fig. 9: **(a) SES convolution** [29,74] The non-trainable basis functions multiply with learnable weights **w** to get kernels. The input then convolves with these kernels to get multi-scale 5D output. **(b) Scale-Projection** [74] takes max over the scale dimension of the 5D output and converts it to 4D. [Key: $*$ = Vanilla convolution.]
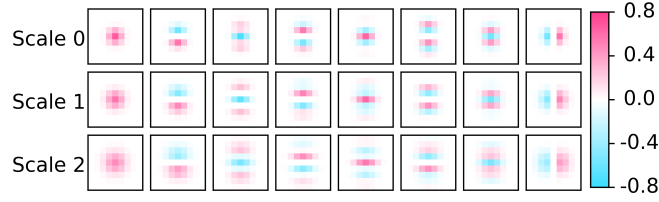


Fig. 10: **Steerable Basis** [74] for $7{\times}7$ SES convolution filters. (Showing only 8 of the 49 members for each scale).

## A2    Implementation Details

We now provide some additional implementation details for facilitating reproduction of this work.

### A2.1    Steerable Filters of SES Convolution

We use the scale equivariant steerable blocks proposed by [73] for our DEVIANT backbone. We now share the implementation details of these steerable filters.
**Basis.** Although steerable filters can use any linearly independent functions as their basis, we stick with the Hermite polynomials as the basis [73]. Let $(0,0)$ denote the center of the function and $(u,v)$ denote the pixel coordinates. Then, the filter coefficients $\psi_{\sigma nm}$ [73] are

$$\psi_{\sigma nm} = \frac{A}{\sigma^2} H_n\left(\frac{u}{\sigma}\right) H_m\left(\frac{v}{\sigma}\right) e^{-\frac{u^2+v^2}{\sigma^2}} \tag{18}$$

$H_n$ denotes the Probabilist's Hermite polynomial of the $n$th order, and $A$ is the normalization constant. The first six Probabilist's Hermite polynomials are

$$H_0(x) = 1 \tag{19}$$

$$H_1(x) = x \tag{20}$$

$$H_2(x) = x^2 - 1 \tag{21}$$

$$H_3(x) = x^3 - 3x \tag{22}$$

$$H_4(x) = x^4 - 6x^2 + 3 \tag{23}$$

Fig. 10 visualizes some of the SES filters and shows that the basis is indeed at different scales.

### A2.2    Monocular 3D Detection

**Architecture.** We use the DLA-34 [98] configuration, with the standard Feature Pyramid Network (FPN) [44], binning and ensemble of uncertainties. FPN is a bottom-up feed-forward CNN that computes feature maps with a downscaling factor of 2, and a top-down network that brings them back to the high-resolution ones. There are total six feature maps levels in this FPN.

   We use DLA-34 as the backbone for our baseline GUP Net [49], while we use SES-DLA-34 as the backbone for DEVIANT. We also replace the 2D pools by 3D pools with pool along the scale dimensions as 1 for DEVIANT.

   We initialize the vanilla CNN from ImageNet weights. For DEVIANT, we use the regularized least squares [73] to initialize the trainable weights in all the Hermite scales from the ImageNet [18] weights. Compared to initializing one of the scales as proposed in [73], we observed more stable convergence in initializing all the Hermite scales.

   We output three foreground classes for KITTI dataset. We also output three foreground classes for Waymo dataset ignoring the Sign class [62].

**Datasets.** We use the publicly available KITTI,Waymo and nuScenes datasets for our experiments. KITTI is available at http://www.cvlibs.net/datasets/kitti/eval_object.php?obj_benchmark=3d under CC BY-NC-SA 3.0 License. Waymo is available at https://waymo.com/intl/en_us/dataset-download-terms/ under the Apache License, Version 2.0. nuScenes is available at https://www.nuscenes.org/nuscenes under CC BY-NC-SA 4.0 International Public License.

**Augmentation.** Unless otherwise stated, we horizontal flip the training images with probability 0.5, and use scale augmentation as 0.4 as well for all the models [49] in training.

**Pre-processing.** The only pre-processing step we use is image resizing.

- *KITTI.* We resize the $[370, 1242]$ sized KITTI images, and bring them to the $[384, 1280]$ resolution [49].
- *Waymo.* We resize the $[1280, 1920]$ sized Waymo images, and bring them to the $[512, 768]$ resolution. This resolution preserves their aspect ratio.

**Box Filtering.** We apply simple hand-crafted rules for filtering out the boxes. We ignore the box if it belongs to a class different from the detection class.

- *KITTI.* We train with boxes which are atleast $2m$ distant from the ego camera, and with visibility $> 0.5$ [49].
- *Waymo.* We train with boxes which are atleast $2m$ distant from the ego camera. The Waymo dataset does not have any occlusion based labels. However,

Waymo provides the number of LiDAR points inside each 3D box which serves as a proxy for the occlusion. We train the boxes which have more than 100 LiDAR points for the vehicle class and have more than 50 LiDAR points for the cyclist and pedestrian class.

**Training.** We use the training protocol of GUP Net [49] for all our experiments. Training uses the Adam optimizer [35] and weight-decay $1 \times 10^{-5}$ . Training dynamically weighs the losses using Hierarchical Task Learning (HTL) [49] strategy keeping $K$ as 5 [49]. Training also uses a linear warmup strategy in the first 5 epochs to stabilize the training. We choose the model saved in the last epoch as our final model for all our experiments.

- *KITTI.* We train with a batch size of 12 on single Nvidia A100 (40GB) GPU for 140 epochs. Training starts with a learning rate $1.25 \times 10^{-3}$ with a step decay of 0.1 at the 90th and the 120th epoch.
- *Waymo.* We train with a batch size of 40 on single Nvidia A100 (40GB) GPU for 30 epochs because of the large size of the Waymo dataset. Training starts with a learning rate $1.25 \times 10^{-3}$ with a step decay of 0.1 at the 18th and the 26th epoch.

**Losses.** We use the GUP Net [49] multi-task losses before the NMS for training. The total loss $\mathcal{L}$ is given by

$$\mathcal{L} = \mathcal{L}_{\text{heatmap}} + \mathcal{L}_{\text{2D,offset}} + \mathcal{L}_{\text{2D,size}} + \mathcal{L}_{\text{3D2D,offset}} + \mathcal{L}_{\text{3D},angle}$$
$$+ \mathcal{L}_{\text{3D},l} + \mathcal{L}_{\text{3D},w} + \mathcal{L}_{\text{3D},h} + \mathcal{L}_{\text{3D},depth}. \tag{24}$$

The individual terms are given by

$$\mathcal{L}_{\text{heatmap}} = \text{Focal}(class^b, class^g), \tag{25}$$

$$\mathcal{L}_{\text{2D,offset}} = L_1(\delta_{\text{2D}}^b, \delta_{\text{2D}}^g), \tag{26}$$

$$\mathcal{L}_{\text{2D,size}} = L_1(w_{\text{2D}}^b, w_{\text{2D}}^g) + L_1(h_{\text{2D}}^b, h_{\text{2D}}^g), \tag{27}$$

$$\mathcal{L}_{\text{3D2D,offset}} = L_1(\delta_{\text{3D2D}}^b, \delta_{\text{3D2D}}^g) \tag{28}$$

$$\mathcal{L}_{\text{3D},angle} = \text{CE}(\alpha^b, \alpha^g) \tag{29}$$

$$\mathcal{L}_{\text{3D},l} = L_1(\mu_{l\text{3D}}^b, \delta_{l\text{3D}}^g) \tag{30}$$

$$\mathcal{L}_{\text{3D},w} = L_1(\mu_{w\text{3D}}^b, \delta_{w\text{3D}}^g) \tag{31}$$

$$\mathcal{L}_{\text{3D},h} = \frac{\sqrt{2}}{\sigma_{h\text{3D}}} L_1(\mu_{h\text{3D}}^b, \delta_{h\text{3D}}^g) + \ln(\sigma_{h\text{3D}}) \tag{32}$$

$$\mathcal{L}_{\text{3D},depth} = \frac{\sqrt{2}}{\sigma_d} L_1(\mu_d^b, \mu_d^g) + \ln(\sigma_d), \tag{33}$$

where,

$$\mu_d^b = f \frac{\mu_{h\text{3D}}^b}{h_{\text{2D}}^b} + \mu_{d,pred} \tag{34}$$

$$\sigma_d = \sqrt{\left(f \frac{\sigma_{h\text{3D}}}{h_{\text{2D}}^b}\right)^2 + \sigma_{d,pred}^2}. \tag{35}$$

The superscripts $b$ and $g$ denote the predicted box and ground truth box respectively. CE and Focal denote the Cross Entropy and Focal loss respectively.

The number of heatmaps depends on the number of output classes. $\delta_{2D}$ denotes the deviation of the 2D center from the center of the heatmap. $\delta_{3D2D,offset}$ denotes the deviation of the projected 3D center from the center of the heatmap. The orientation loss is the cross entropy loss between the binned observation angle of the prediction and the ground truth. The observation angle $\alpha$ is split into 12 bins covering $30°$ range. $\delta_{l3D}, \delta_{w3D}$ and $\delta_{h3D}$ denote the deviation of the 3D length, width and height of the box from the class dependent mean size respectively.

The depth is the hardest parameter to estimate [53]. So, GUP Net uses in-network ensembles to predict the depth. It obtains a Laplacian estimate of depth from the 2D height, while it obtains another estimate of depth from the prediction of depth. It then adds these two depth estimates.

**Inference.** Our testing resolution is same as the training resolution. We do not use any augmentation for test/validation. We keep the maximum number of objects to 50 in an image, and we multiply the class and predicted confidence to get the box's overall score in inference as in [36]. We consider output boxes with scores greater than a threshold of 0.2 for KITTI [49] and 0.1 for Waymo [62].

Table 14: **Generalization gap** (↓) on KITTI Val cars. Monocular detection has huge generalization gap between training and inference sets. [Key: **Best**]

| Method | Scale Eqv | Set | IoU$_{3D}$ ≥ 0.7 | | | | | | IoU$_{3D}$ ≥ 0.5 | | | | | |
| | | | AP$_{3D|R_{40}}$[%](↑) | | | AP$_{BEV|R_{40}}$[%](↑) | | | AP$_{3D|R_{40}}$[%](↑) | | | AP$_{BEV|R_{40}}$[%](↑) | | |
| | | | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard |
| GUP Net [49] | | Train | 91.83 | 74.87 | 67.43 | 95.19 | 80.95 | 73.55 | 99.50 | 93.62 | 86.22 | 99.56 | 93.88 | 86.46 |
| | | Val | 21.10 | 15.48 | 12.88 | 28.58 | 20.92 | 17.83 | 58.95 | 43.99 | 38.07 | 64.60 | 47.76 | 42.97 |
| | | Gap | 70.73 | **59.39** | 54.55 | 66.61 | 60.03 | 55.72 | 40.55 | 49.63 | **48.15** | 34.96 | 46.12 | **43.49** |
| **DEVIANT** | ✓ | Train | 91.09 | 76.19 | 67.16 | 94.76 | 82.61 | 75.51 | 99.37 | 93.56 | 88.57 | 99.50 | 93.87 | 88.90 |
| | | Val | 24.63 | 16.54 | 14.52 | 32.60 | 23.04 | 19.99 | 61.00 | 46.00 | 40.18 | 65.28 | 49.63 | 43.50 |
| | | Gap | **66.46** | 59.65 | **52.64** | **62.16** | **59.57** | **55.52** | **38.37** | **47.56** | 48.39 | **34.22** | **44.24** | 45.40 |

Table 15: **Comparison on multiple backbones** on KITTI Val cars. [Key: **Best**]

| BackBone | Method | IoU$_{3D}$ ≥ 0.7 | | | | | | IoU$_{3D}$ ≥ 0.5 | | | | | |
| | | AP$_{3D|R_{40}}$[%](↑) | | | AP$_{BEV|R_{40}}$[%](↑) | | | AP$_{3D|R_{40}}$[%](↑) | | | AP$_{BEV|R_{40}}$[%](↑) | | |
| | | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard |
| ResNet-18 | GUP Net [49] | 18.86 | 13.20 | 11.01 | 26.05 | 19.37 | 16.57 | 54.90 | 40.65 | 34.98 | 60.54 | 46.13 | 40.12 |
| | **DEVIANT** | 20.27 | 14.21 | 12.56 | 28.09 | 20.32 | 17.49 | 55.75 | 42.41 | 36.97 | 60.82 | 46.43 | 40.59 |
| DLA-34 | GUP Net [49] | 21.10 | 15.48 | 12.88 | 28.58 | 20.92 | 17.83 | 58.95 | 43.99 | 38.07 | 64.60 | 47.76 | 42.97 |
| | **DEVIANT** | **24.63** | **16.54** | **14.52** | **32.60** | **23.04** | **19.99** | **61.00** | **46.00** | **40.18** | **65.28** | **49.63** | **43.50** |

## A3    Additional Experiments and Results

We now provide additional details and results of the experiments evaluating our system's performance.

### A3.1    KITTI Val Split

**Monocular Detection has Huge Generalization Gap.** As mentioned in Sec. 1, we now show that the monocular detection has huge generalization gap between training and inference. We report the object detection performance on the train and validation (val) set for the two models on KITTI Val split in Tab. 14. Tab. 14 shows that the performance of our baseline GUP Net [49] and our DEVIANT is huge on the training set, while it is less than one-fourth of the train performance on the val set.

We also report the generalization gap metric [93] in Tab. 14, which is the difference between training and validation performance. The generalization gap at both the thresholds of 0.7 and 0.5 is huge.

**Comparison on Multiple Backbones.** A common trend in 2D object detection community is to show improvements on multiple backbones [82]. DD3D [57] follows this trend and also reports their numbers on multiple backbones. Therefore, we follow the same and compare with our baseline on multiple backbones on KITTI Val cars in Tab. 15. Tab. 15 shows that DEVIANT shows consistent improvements over GUP Net [49] in 3D object detection on multiple backbones, proving the effectiveness of our proposal.

**Comparison with Bigger CNN Backbones.** Since the SES blocks increase the Flop counts significantly compared to the vanilla convolution block, we next compare DEVIANT with bigger CNN backbones with comparable GFLOPs and FPS/ wall-clock time (instead of same configuration) in Tab. 16. We compare

Table 16: **Results with bigger CNNs having similar flops** on KITTI Val cars. [Key: **Best**]

| Method | BackBone | Param (↓) (M) | Disk Size (↓) (MB) | Flops (↓) (G) | Infer (↓) (ms) | AP$_{3D}$ IoU$_{3D} \geq$ 0.7 (↑) Easy | Mod | Hard | AP$_{3D}$ IoU$_{3D} \geq$ 0.5 (↑) Easy | Mod | Hard |
|---|---|---|---|---|---|---|---|---|---|---|---|
| GUP Net [49] | DLA-34 | **16** | **235** | **30** | **20** | 21.10 | 15.48 | 12.88 | 58.95 | 43.99 | 38.07 |
| GUP Net [49] | DLA-102 | 34 | 583 | 70 | 25 | 20.96 | 14.64 | 12.80 | 57.06 | 41.78 | 37.26 |
| GUP Net [49] | DLA-169 | 54 | 814 | 114 | 30 | 21.76 | 15.35 | 12.72 | 57.60 | 43.27 | 37.32 |
| DEVIANT | SES-DLA-34 | **16** | 236 | 235 | 40 | **24.63** | **16.54** | **14.52** | **61.00** | **46.00** | **40.18** |

Table 17: **Results on KITTI Val cyclists and pedestrians** (Cyc/Ped) (IoU$_{3D} \geq$ 0.5). [Key: **Best**, **Second Best**]

| Method | Extra | Cyc AP$_{3D|R_{40}}$[%](↑) Easy | Mod | Hard | Ped AP$_{3D|R_{40}}$[%](↑) Easy | Mod | Hard |
|---|---|---|---|---|---|---|---|
| GrooMeD-NMS [36] | − | 0.00 | 0.00 | 0.00 | 3.79 | 2.71 | 2.61 |
| MonoDIS [70] | − | 1.52 | 0.73 | 0.71 | 3.20 | 2.28 | 1.71 |
| MonoDIS-M [68] | − | 2.70 | 1.50 | 1.30 | 9.50 | 7.10 | 5.70 |
| GUP Net (Retrained) [49] | − | **4.41** | **2.17** | **2.03** | **9.37** | **6.84** | **5.73** |
| **DEVIANT (Ours)** | − | **4.05** | **2.20** | **2.14** | **9.85** | **7.18** | **5.42** |

DEVIANT with DLA-102 and DLA-169 - two biggest DLA networks with ImageNet weights[4] on KITTI Val split. We use the fvcore library[5] to get the parameters and flops. Tab. 16 shows that DEVIANT again outperforms the bigger CNN backbones, especially on nearby objects. We believe this happens because the bigger CNN backbones have more trainable parameters than DEVIANT, which leads to overfitting. Although DEVIANT takes more time compared to the CNN backbones, DEVIANT still keeps the inference almost real-time.

**Performance on Cyclists and Pedestrians.** Tab. 17 lists out the results of 3D object detection on KITTI Val Cyclist and Pedestrians. The results show that DEVIANT is competitive on challenging Cyclist and achieves SOTA results on Pedestrians on the KITTI Val split.

**Cross-Dataset Evaluation Details.** For cross-dataset evaluation, we test on all 3,769 images of the KITTI Val split, as well as all frontal 6,019 images of the nuScenes Val split [9], as in [67]. We first convert the nuScenes Val images to the KITTI format using the `export_kitti`[6] function in the nuscenes devkit. We keep KITTI Val images in the $[384, 1280]$ resolution, while we keep the nuScenes Val images in the $[384, 672]$ resolution to preserve the aspect ratio. For M3D-RPN [4], we bring the nuScenes Val images in the $[512, 910]$ resolution.

Monocular 3D object detection relies on the camera focal length to back-project the projected centers into the 3D space. Therefore, the 3D centers depends on the focal length of the camera used in the dataset. Hence, one should take the camera focal length into account while doing cross-dataset evaluation. We now calculate the camera focal length of a dataset as follows. We take the camera matrix **K** and calculate the normalized focal length $\bar{f} = \frac{2f_y}{H}$, where $H$ denotes the height of the image. The normalized focal length $\bar{f}$ for the KITTI

---

[4] Available at http://dl.yf.io/dla/models/imagenet/

[5] https://github.com/facebookresearch/fvcore

[6] https://github.com/nutonomy/nuscenes-devkit/blob/master/python-sdk/nuscenes/scripts/export_kitti.py

Table 18: **Stress Test** with rotational and $xy$-translation ego movement on KITTI Val cars. [Key: **Best**]

| Set | Method | AP$_{3D}$ IoU$_{3D}\geq 0.7$ (↟) | | | AP$_{3D}$ IoU$_{3D}\geq 0.5$ (↟) | | |
|-----|--------|------|------|------|------|------|------|
| | | Easy | Mod | Hard | Easy | Mod | Hard |
| Subset | GUP Net [49] | 17.22 | 11.43 | 9.91 | 47.47 | 35.02 | 32.63 |
| (306) | DEVIANT | **20.17** | **12.49** | **10.93** | **49.81** | **36.93** | **34.32** |
| KITTI Val | GUP Net [49] | 21.10 | 15.48 | 12.88 | 58.95 | 43.99 | 38.07 |
| (3769) | DEVIANT | **24.63** | **16.54** | **14.52** | **61.00** | **46.00** | **40.18** |

dataset is 3.82, while the normalized focal length $\bar{f}$ for the nuScenes dataset is 2.82. Thus, the KITTI and the nuScenes images have a different focal length [84].

M3D-RPN [4] does not normalize w.r.t. the focal length. So, we explicitly correct and divide the depth predictions of nuScenes images from the KITTI model by $3.82/2.82 = 1.361$ in the M3D-RPN [4] codebase. The GUP Net [49] and DEVIANT codebases use normalized coordinates *i.e.* they normalize w.r.t. the focal length. So, we do not explicitly correct the focal length for GUP Net and DEVIANT predictions.

We match predictions to the ground truths using the IoU$_{2D}$ overlap threshold of 0.7 [67]. After this matching, we calculate the Mean Average Error (MAE) of the depths of the predicted and the ground truth boxes [67].

**Stress Test with Rotational and/or xy-translation Ego Movement.** Corollary 1 uses translation along the depth as the sole ego movement. This assumption might be valid for the current outdoor datasets and benchmarks, but is not the case in the real world. Therefore, we conduct stress tests on how tolerable DEVIANT and GUP Net [49] are when there is rotational and/or $xy$-translation movement on the vehicle.

First, note that KITTI and Waymo are already large-scale real-world datasets, and our own dataset might not be a good choice. So, we stick with KITTI and Waymo datasets. We manually choose 306 KITTI Val images with such ego movements and again compare performance of DEVIANT and GUP Net on this subset in Tab. 18. The average distance of the car in this subset is 27.69 m ($\pm$16.59 m), which suggests a good variance and unbiasedness in the subset. Tab. 18 shows that both the DEVIANT backbone and the CNN backbone show a drop in the detection performance by about 4 AP points on the Mod cars of ego-rotated subset compared to the all set. This drop experimentally confirms the theory that both the DEVIANT backbone and the CNN backbone do not handle arbitrary 3D rotations. More importantly, the table shows that DEVIANT maintains the performance improvement over GUP Net [49] under such movements.

Also, Waymo has many images in which the ego camera shakes. Improvements on Waymo (Tab. 12) also confirms that DEVIANT outperforms GUP Net [49] even when there is rotational or $xy$-translation ego movement.

**Comparison of Depth Estimates from Monocular Depth Estimators and 3D Object Detectors.** We next compare the depth estimates from monocular depth estimators and depth estimates from monocular 3D object detectors on the foreground objects. We take a monocular depth estimator BTS [41] model trained on KITTI Eigen split. We next compare the depth error for all and fore-

Table 19: **Comparison of Depth Estimates** of monocular depth estimators and 3D object detectors on KITTI Val cars. Depth from a depth estimator BTS is not good for foreground objects (cars) beyond 20+ m range. [Key: **Best**, **Second Best**]

| Method | Depth at | Ground Truth | Back+ Foreground | | | Foreground (Cars) | | |
|---|---|---|---|---|---|---|---|---|
| | | | $0-20$ | $20-40$ | $40-\infty$ | $0-20$ | $20-40$ | $40-\infty$ |
| GUP Net [49] | 3D Center | 3D Box | – | – | – | 0.45 | **1.10** | **1.85** |
| DEVIANT | 3D Center | 3D Box | – | – | – | **0.40** | **1.09** | **1.80** |
| BTS [41] | Pixel | LiDAR | 0.48 | 1.30 | 1.83 | **0.30** | 1.22 | 2.16 |



Fig. 11: **Equivariance error** ($\Delta$) comparison for DEVIANT and GUP Net on previous three frames of the KITTI monocular videos at block 3 in the backbone.

ground objects (cars) on KITTI Val split using MAE (↓) metric in Tab. 19 as in Tab. 6. We use the MSeg [39] to segment out cars in the driving scenes for BTS. Tab. 19 shows that the depth from BTS is not good for foreground objects (cars) beyond 20+ m range. Note that there is a data leakage issue between the KITTI Eigen train split and the KITTI Val split [69] and therefore, we expect more degradation in performance of monocular depth estimators after fixing the data leakage issue.

**Equivariance Error for KITTI Monocular Videos.** A better way to compare the scale equivariance of the DEVIANT and GUP Net [49] compared to Fig. 4, is to compare equivariance error on real images with depth translations of the ego camera. The equivariance error $\Delta$ is the normalized difference between the scaled feature map and the feature map of the scaled image, and is given by

$$\Delta = \frac{1}{N} \sum_{i=1}^{N} \frac{||\mathcal{T}_{s_i}\Phi(h_i) - \Phi(\mathcal{T}_{s_i}h_i)||_2^2}{||\mathcal{T}_{s_i}\Phi(h_i)||_2^2}, \tag{36}$$

where $\Phi$ denotes the neural network, $\mathcal{T}_{s_i}$ is the scaling transformation for the image $i$, and $N$ is the total number of images. Although we do evaluate this error in Fig. 4, the image scaling in Fig. 4 does not involve scene change because of the absence of the moving objects. Therefore, evaluating on actual depth translations of the ego camera makes the equivariance error evaluation more realistic. We next carry out this experiment and report the equivariance error on three previous frames of the val images of the KITTI Val split as in [5]. We plot this equivariance error in Fig. 11 at block 3 of the backbones because the resolution at this block corresponds to the output feature map of size $[96, 320]$. Fig. 11 is similar to Fig. 4b, and shows that DEVIANT achieves lower equivariance error. Therefore,

Table 20: **Five Different Runs** on KITTI Val cars. [Key: **Average**]

| Method | Run | IoU$_{3D} \geq 0.7$ | | | | | | IoU$_{3D} \geq 0.5$ | | | | | |
| | | AP$_{3D|R_{40}}$[%](↑) | | | AP$_{BEV|R_{40}}$[%](↑) | | | AP$_{3D|R_{40}}$[%](↑) | | | AP$_{BEV|R_{40}}$[%](↑) | | |
| | | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard | Easy | Mod | Hard |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GUP Net [49] | 1 | 21.67 | 14.75 | 12.68 | 28.72 | 20.88 | 17.79 | 58.27 | 43.53 | 37.62 | 63.67 | 47.37 | 42.55 |
| | 2 | 21.26 | 14.94 | 12.49 | 28.39 | 20.40 | 17.43 | 59.20 | 43.55 | 37.63 | 64.06 | 47.46 | 42.67 |
| | 3 | 20.87 | 15.03 | 12.61 | 28.66 | 20.56 | 17.48 | 60.19 | 44.08 | 39.36 | 65.26 | 49.44 | 43.17 |
| | 4 | 21.10 | 15.48 | 12.88 | 28.58 | 20.92 | 17.83 | 58.95 | 43.99 | 38.07 | 64.60 | 47.76 | 42.97 |
| | 5 | 22.52 | 15.92 | 13.31 | 30.77 | 22.40 | 19.36 | 59.91 | 44.00 | 39.30 | 64.94 | 48.01 | 43.08 |
| | Avg | **21.48** | **15.22** | **12.79** | **29.02** | **21.03** | **17.98** | **59.30** | **43.83** | **38.40** | **64.51** | **48.01** | **42.89** |
| DEVIANT | 1 | 23.19 | 15.84 | 14.11 | 29.82 | 21.93 | 19.16 | 60.19 | 45.52 | 39.86 | 66.32 | 49.39 | 43.38 |
| | 2 | 23.33 | 16.12 | 13.54 | 31.22 | 22.64 | 19.64 | 61.59 | 46.33 | 40.35 | 67.49 | 50.26 | 43.98 |
| | 3 | 24.12 | 16.37 | 14.48 | 31.58 | 22.52 | 19.65 | 62.51 | 46.47 | 40.65 | 67.33 | 50.24 | 44.16 |
| | 4 | 24.63 | 16.54 | 14.52 | 32.60 | 23.04 | 19.99 | 61.00 | 46.00 | 40.18 | 65.28 | 49.63 | 43.50 |
| | 5 | 25.82 | 17.69 | 15.07 | 33.63 | 23.84 | 20.60 | 62.39 | 46.46 | 40.61 | 67.55 | 50.51 | 45.80 |
| | Avg | **24.22** | **16.51** | **14.34** | **31.77** | **22.79** | **19.81** | **61.54** | **46.16** | **40.33** | **66.79** | **50.01** | **44.16** |

Table 21: **Experiments Comparison**.

| Method | Venue | Multi-Dataset | Cross-Dataset | Multi-Backbone |
|---|---|---|---|---|
| GrooMeD-NMS [36] | CVPR21 | − | − | − |
| MonoFlex [100] | CVPR21 | − | − | − |
| CaDDN [62] | CVPR21 | ✓ | − | − |
| MonoRCNN [67] | ICCV21 | − | ✓ | − |
| GUP Net [49] | ICCV21 | − | − | − |
| DD3D [57] | ICCV21 | ✓ | − | ✓ |
| PCT [80] | NeurIPS21 | ✓ | − | ✓ |
| MonoDistill [14] | ICLR22 | − | − | − |
| MonoDIS-M [68] | TPAMI20 | ✓ | − | − |
| MonoEF [103] | TPAMI21 | ✓ | − | − |
| **DEVIANT** | - | ✓ | ✓ | ✓ |

DEVIANT has better equivariance to depth translations (scale transformation s) than GUP Net [49] in real scenarios.

**Model Size, Training, and Inference Times.** Both DEVIANT and the baseline GUP Net have the same number of trainable parameters, and therefore, the same model size. GUP Net takes 4 hours to train on KITTI Val and 0.02 ms per image for inference on a single Ampere A100 (40 GB) GPU. DEVIANT takes 8.5 hours for training and 0.04 ms per image for inference on the same GPU. This is expected because SE models use more flops [74, 104] and, therefore, DEVIANT takes roughly twice the training and inference time as GUP Net.

**Reproducibility.** As described in Sec. 5.2, we now list out the five runs of our baseline GUP Net [49] and DEVIANT in Tab. 20. Tab. 20 shows that DEVIANT outperforms GUP Net in all runs and in the average run.

**Experiment Comparison.** We now compare the experiments of different papers in Tab. 21. To the best of our knowledge, the experimentation in DEVIANT is more than the experimentation of most monocular 3D object detection papers.

## A3.2    Qualitative Results

**KITTI.** We next show some more qualitative results of models trained on KITTI Val split in Fig. 13. We depict the predictions of DEVIANT in image view on the left and the predictions of DEVIANT and GUP Net [49], and ground truth
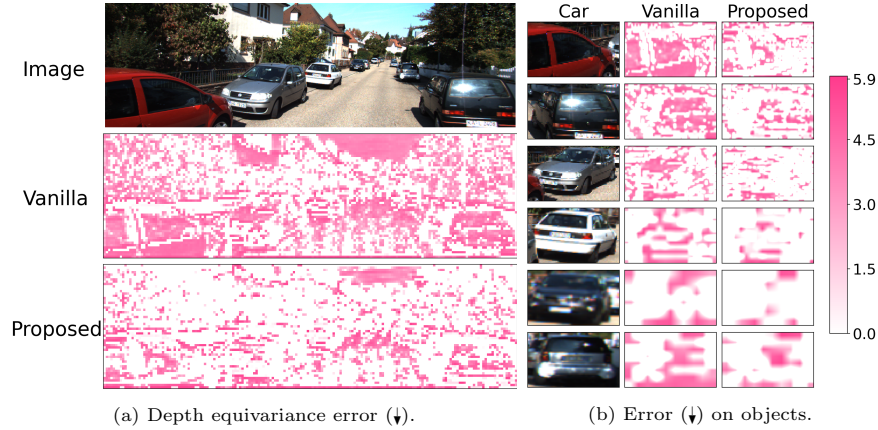
(a) Depth equivariance error (↓).          (b) Error (↓) on objects.

Fig. 12: **(a) Depth (scale) equivariance error** of vanilla GUP Net [49] and proposed DEVIANT. (See Sec. 5.2 for details) **(b) Error on objects.** The proposed backbone has less depth equivariance error than vanilla CNN backbone.

in BEV on the right. In general, DEVIANT predictions are more closer to the ground truth than GUP Net [49].

**nuScenes Cross-Dataset Evaluation.** We then show some qualitative results of KITTI Val model evaluated on nuScenes frontal in Fig. 14. We again observe that DEVIANT predictions are more closer to the ground truth than GUP Net [49]. Also, considerably less number of boxes are detected in the cross-dataset evaluation *i.e.* on nuScenes. We believe this happens because of the domain shift.

**Waymo.** We now show some qualitative results of models trained on Waymo Val split in Fig. 15. We again observe that DEVIANT predictions are more closer to the ground truth than GUP Net [49].

### A3.3    Demo Videos of DEVIANT

**Detection Demo.** We next put a short demo video of our DEVIANT model trained on KITTI Val split at https://www.youtube.com/watch?v=2D73ZBrU-PA. We run our trained model independently on each frame of 2011_09_26_drive_0009 KITTI raw [27]. The video belongs to the City category of the KITTI raw video. None of the frames from the raw video appear in the training set of KITTI Val split [36]. We use the camera matrices available with the video but do not use any temporal information. Overlaid on each frame of the raw input videos, we plot the projected 3D boxes of the predictions and also plot these 3D boxes in the BEV. We set the frame rate of this demo at 10 fps as in KITTI. The attached demo video demonstrates very stable and impressive results because of the additional equivariance to depth translations in DEVIANT which is absent in vanilla CNNs. Also, notice that the orientation of the boxes are stable despite not using any temporal information.

**Equivariance Error Demo.** We next show the depth equivariance (scale equivariance) error demo of one of the channels from the vanilla GUP Net and our

proposed method at https://www.youtube.com/watch?v=7ODIjQkuZvw. As before, we report at block 3 of the backbones which corresponds to output feature map of the size [96, 320]. The equivariance error demo indicates more white spaces which confirms that DEVIANT achieves lower equivariance error compared to the baseline GUP Net [49]. Thus, this demo agrees with Fig. 12a. This happens because depth (scale) equivariance is additionally hard-baked into DEVIANT, while the vanilla GUP Net is not equivariant to depth translations (scale transformations).

## Acknowledgements

Fig. 13: **KITTI Qualitative Results**. DEVIANT predictions in general are more accurate than GUP Net [49]. [Key: Cars, Cyclists and Pedestrians of DEVIANT; all classes of GUP Net, and Ground Truth in BEV].
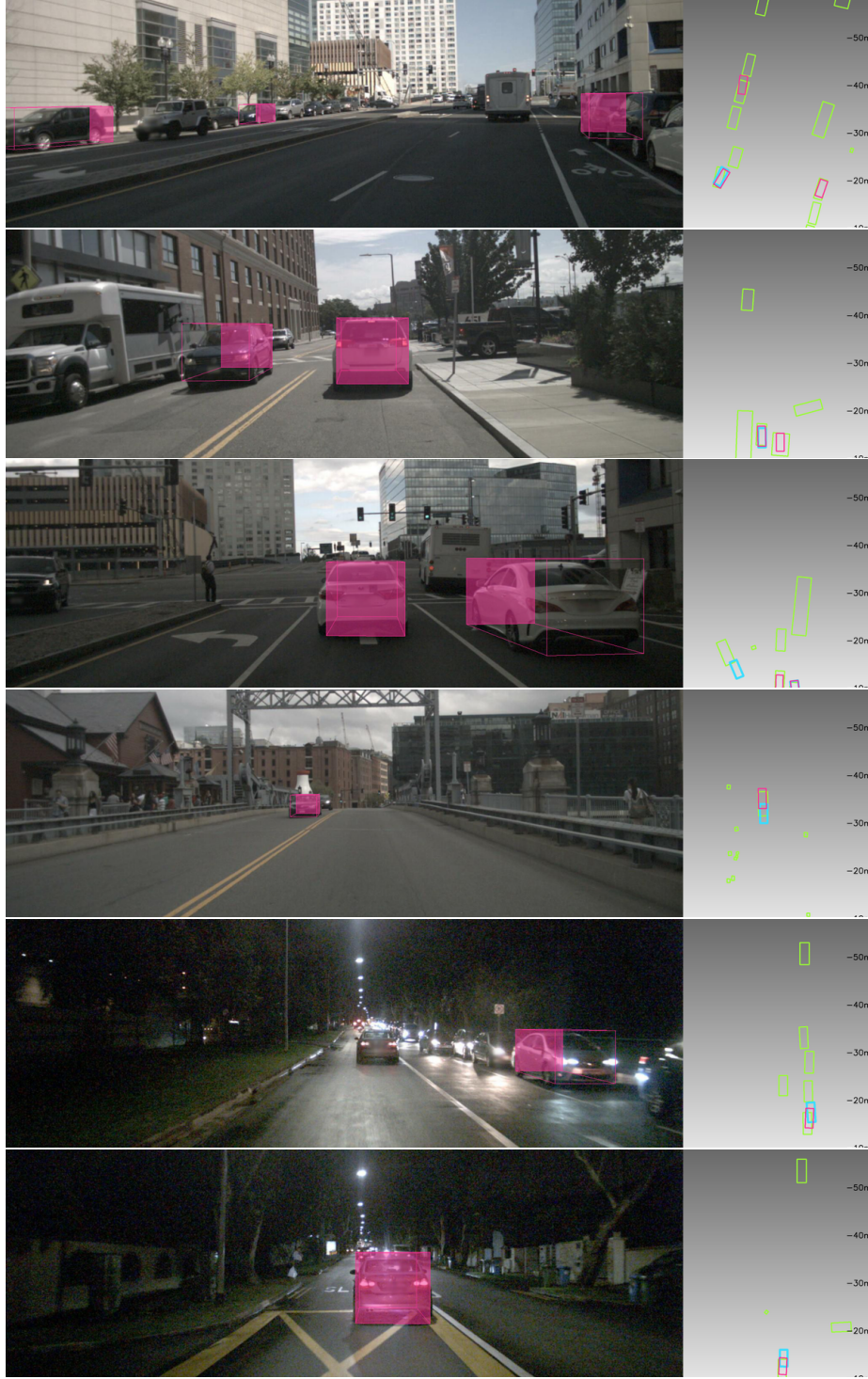
Fig. 14: **nuScenes Cross-Dataset Qualitative Results.** DEVIANT predictions in general are more accurate than GUP Net [49]. [Key: Cars of DEVIANT; Cars of GUP Net, and Ground Truth in BEV].
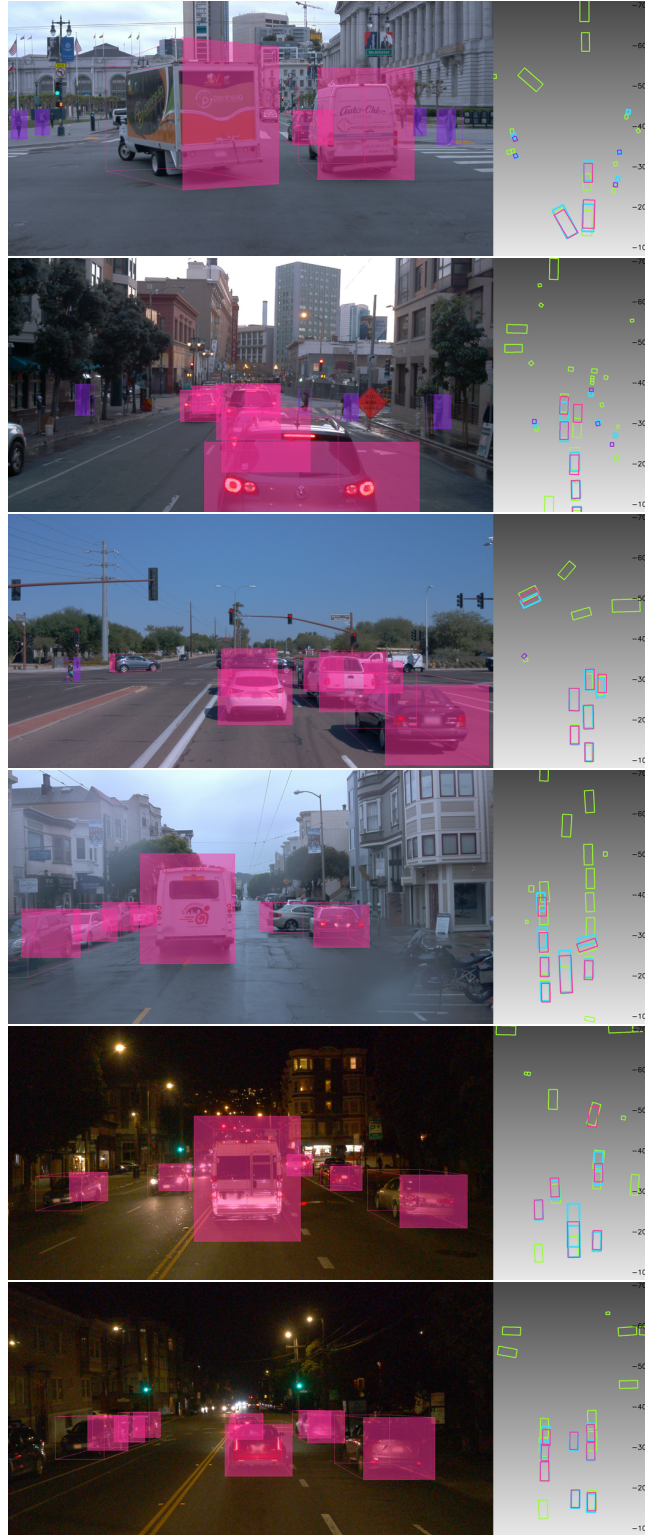
Fig. 15: **Waymo Qualitative Results**. DEVIANT predictions in general are more accurate than GUP Net [49]. [Key: Cars, Cyclists and Pedestrians of DEVIANT; all classes of GUP Net, and Ground Truth in BEV].