# GBlobs: Explicit Local Structure via Gaussian Blobs for Improved Cross-Domain LiDAR-based 3D Object Detection

Dušan Malić[1,2]      Christian Fruhwirth-Reisinger[1,2]      Samuel Schulter[3,†]      Horst Possegger[1,2]

[1]Christian Doppler Laboratory for Embedded Machine Learning
[2]Institute of Visual Computing, Graz University of Technology
[3]Amazon

{dusan.malic, reisinger, possegger}@tugraz.at

## Abstract

*LiDAR-based 3D detectors need large datasets for training, yet they struggle to generalize to novel domains. Domain Generalization (DG) aims to mitigate this by training detectors that are invariant to such domain shifts. Current DG approaches exclusively rely on **global** geometric features (point cloud Cartesian coordinates) as input features. Over-reliance on these global geometric features can, however, cause 3D detectors to prioritize object location and absolute position, resulting in poor cross-domain performance. To mitigate this, we propose to exploit explicit **local** point cloud structure for DG, in particular by encoding point cloud neighborhoods with Gaussian blobs, **GBlobs**. Our proposed formulation is highly efficient and requires no additional parameters. Without any bells and whistles, simply by integrating **GBlobs** in existing detectors, we beat the current state-of-the-art in challenging single-source DG benchmarks by over 21 mAP (Waymo→KITTI), 13 mAP (KITTI→Waymo), and 12 mAP (nuScenes→KITTI), without sacrificing in-domain performance. Additionally, GBlobs demonstrate exceptional performance in multi-source DG, surpassing the current state-of-the-art by 17, 12, and 5 mAP on Waymo, KITTI, and ONCE, respectively.*

## 1. Introduction

LiDAR-based 3D detection models [4, 46, 48, 61, 65] predominantly rely on *global* input features, *i.e.* representing points by their Cartesian coordinates. Only a small subset of research *additionally* incorporates *local* point cloud information, like relative distance to the voxel center [8, 55], LiDAR intensity [38], surface normals [37, 69] or handcrafted local descriptors [25] to enhance model performance. The convenience of global input features and their strong

(a) Our proposed encoder estimates Gaussian blobs from local neighborhoods, providing a more expressive representation of local geometry compared to PointNet [40] and the more commonly used mean encoding.



(b) Comparison with state-of-the-art single-source DG methods. Benchmark results for Waymo→KITTI (W→K), KITTI→Waymo (K→W), and nuScenes→KITTI (n→K).

(c) Comparison with state-of-the-art multi-sorce DG methods. Benchmark results from training on all (∗), excluding the target dataset, and evaluating it on ONCE (O), KITTI (K) and Waymo (W).
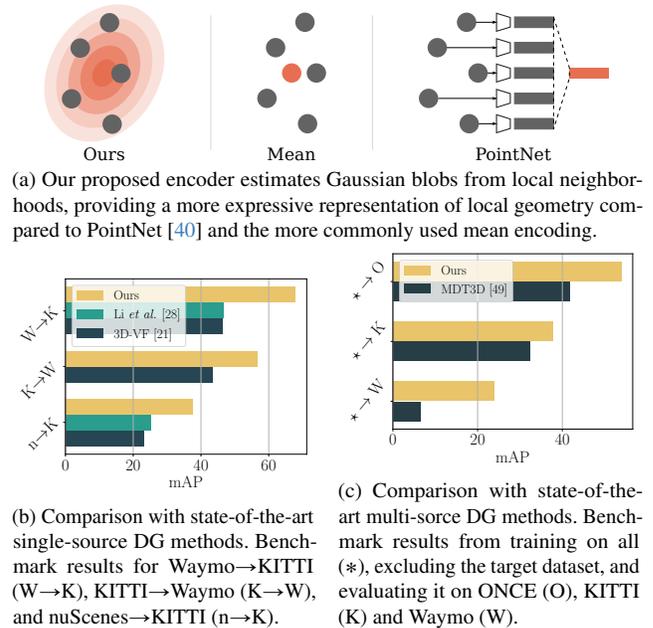
Figure 1. We propose a novel encoding scheme for local neighborhood in point clouds, using a Gaussian blob instead of the commonly employed mean or PointNet encoding (a). Our simple yet effective approach significantly outperforms state-of-the-art methods in both single- and multi-source domain generalization does not require any additional parameters, spatial training procedures, or hyperparameter optimization (b, c).

in-domain performance established this representation as the default choice for most state-of-the-art detection models. However, such models are highly sensitive to rigid transformations [25, 57], being more biased toward object position than local features like shape or appearance. Therefore, they generalize poorly to novel, unseen domains [9, 56, 62, 63].

In this paper, we demonstrate that local point cloud geometry can significantly improve model generalization. In particular, we introduce a new method for representing lo-

cal neighborhoods as Gaussian blobs (see Fig. 1a), defined by their means and covariance matrices. Our representation effectively decouples object location from model encoding, enabling the model to focus on learning local shape and appearance. Unlike existing local descriptors based on relative distances [8, 55] or surface normals [15], our formulation is more descriptive due to the covariance modeling and eliminates ambiguities arising from surface orientation. Moreover, our method is inherently invariant to point permutations (*i.e.* a neighborhood's covariance matrix remains the same regardless of the point order), which eliminates the need for point cloud sorting [25, 26, 43] or channel-wise max pooling and concatenation [40] (in Fig. 1a). To demonstrate the effectiveness of our approach, we evaluate it on challenging cross-domain LiDAR 3D object detection benchmarks, highlighting the significant improvements on *domain generalization* (DG) for a wide range of 3D detectors. While previous methods often require specialized data augmentation [21, 49], auxiliary tasks [28], or domain-invariant feature learning [59], which can introduce additional complexity and limitations, our module is simple, effective, universally applicable to any dataset and model, and does not degrade the in-domain performance.

In the single-domain generalization benchmarks, as shown in Fig. 1b, we achieve significant improvements over the state-of-the-art, with over 21 mAP in Waymo→KITTI (W→K), 13 mAP in KITTI→Waymo (K→W), and 12 mAP in nuScenes→KITTI (n→K). In our multi-domain DG experiments in Fig. 1c, where we trained on all datasets (∗) except the target dataset, we consistently observed significant improvements over the state-of-the-art, with gains of over 17, 12, and 5 mAP on Waymo, KITTI, and ONCE, respectively. In summary, our main contributions are:

- A novel, permutation-invariant LiDAR point cloud representation using Gaussian blobs, offering descriptive features, efficient computation, and no additional model parameters.
- Comprehensive evaluations on both single- and multisource domain generalization benchmarks and in-depth ablation studies.
- Publicly accessible code and trained models at https://github.com/malicd/GBlobs

## 2. Related Work

**LiDAR-based 3D object detection** can be broadly classified by their input representations. Methods such as [41, 46, 60, 64] operate directly on unordered point cloud data. To address the computational challenges of processing raw point clouds, [36] introduces a LiDAR range view representation and leverages efficient 2D convolutions. Similarly, pillar-based networks [10, 20, 24, 51] encode 3D data into a sparse 2D grid where each grid cell (pillar) spans the entire height axis. To address the information loss inherent in such

dimension reduction, voxel-based methods [4, 8, 55, 61, 65] directly process 3D voxels derived from input point clouds. 3D object detectors commonly use Cartesian coordinates as input features.

There exists a handful of studies that, in addition to utilizing Cartesian coordinates, also incorporates local point cloud geometry. VoxelNet [70] augments the input data by appending the relative distance between each point and its voxel's centroid to its global coordinates. Instead of voxel's centroid, PointPillars [20] utilizes distance to the voxel center. To capture more local information, several methods [20, 55, 70] employ a shared PointNet [40] architecture to encode pillars or voxels. [37, 69] utilize a two stream approach, encoding Cartesian coordinates and surface normals in separate branches and fusing them in later stages. Hybrid methods, which combine local and global geometric inputs, share a common limitation with methods trained solely on global features: they often exhibit poor generalization across diverse domains, likely because their performance is heavily influenced by global features.

In contrast, our proposed approach is *exclusively* based on local point cloud geometry, demonstrating the crucial role of these features in model generalization. We show that our representation provides more detail (than, *e.g.*, simple relative distance) and is less susceptible to noise (which degrades performance of surface normal-based features).

**Unsupervised Domain Adaptation (UDA)** transfers knowledge from a labeled source domain to an unlabeled target domain, using both domains to adapt models to the novel data distributions. Substantial efforts have focused on explicitly addressing different challenges, *e.g.* object size bias [34, 56, 56, 62] or overcoming varying LiDAR resolutions [16, 44, 58]. In UDA, self-training remains the dominant paradigm: This typically involves a closed-loop process where the pseudo-label database is iteratively updated with predictions generated by a model trained on these pseudolabels [3, 5, 62, 63, 68]. To enhance pseudo-label quality, several methods exploit the temporal dimension inherent in autonomous driving data [11, 54, 66].

Despite a comprehensive portfolio of LiDAR-based UDA approaches, the impact of using different input features on model performance remains largely unexplored. Most existing methods rely solely on the standard global point cloud inputs, neglecting other feature channels. Given the prevalence of self-training in UDA, good domain generalization is crucial to obtain improved pseudo labels. We demonstrate how our proposed local feature representation significantly enhances the domain generalization of various models, which could be leveraged for improved UDA via common selftraining loops.

**Domain Generalization (DG)** aims to train models that can effectively generalize to new, unseen domains without requiring any additional data from those domains. For 2D image-

based tasks, DG research primarily focuses on learning domain-invariant features [13, 23, 29, 42] or disentangling domain-specific and domain-invariant features [32, 33, 39]. More recently, studies such as [1, 6] have sought to enhance generalization by employing powerful vision-language models. A smaller subset of DG research explores data augmentation techniques to improve generalization [17, 53, 67].

In the context of LiDAR-based detection, data augmentation [14, 22] is crucial to prevent model overfitting to training data. To this end, 3D-VF [21] employs adversarial augmentation to generate diverse and realistic training examples from a single source domain. In addition to a novel augmentation strategy, Li *et al.* [28] enhance the model robustness by incorporating an auxiliary adversarial task during training and at test-time. Similarly, MDT3D [49] leverages multiple source domains and cross-dataset augmentation to learn more robust models. Wu *et al.* [59] also utilize multiple source domains to learn more general, dataset-invariant features.

A shared characteristic of DG methods for LiDAR-based object detection is their reliance solely on global point coordinates as network inputs. While the intensity channel has been shown to have a negative impact on domain generalization [18, 49], the influence of local geometric features has not been explored so far. Our simple yet effective local feature representation significantly enhances model generalization without requiring tailored data augmentation, auxiliary training tasks, nor test-time adaptation, demonstrating that a focus on local geometric information can be more beneficial than complex DG strategies.

## 3. Explicit Local Structure for Better DG

LiDAR Domain Generalization (DG) is a critical research area for improving the robustness of detection models across domains. Prior work often relies solely on global point cloud features (*i.e.* Cartesian coordinates) and techniques like data augmentation [21, 49] or auxiliary tasks [28]. In contrast, we propose leveraging local point cloud geometry to improve the generalization capabilities.

Local features force models to learn object shapes independently of their global position, thereby reducing sensitivity to rigid transformations. This shape-based representation enhances the model's ability to generalize to new domains. However, common local features, such as relative distance [8] or RepSurf [43], are often used in conjunction with the global representation, hindering their domain generalization effect. While highly efficient, relative distance features are typically outperformed by more sophisticated representation of local geometry, such as RepSurf (as shown in our experiments). However, RepSurf leverages surface normal estimates, which are sensitive to noise (due to ambiguity in surface orientation [15, 27]), and requires point cloud sorting, causing it to be a computationally expensive input representation.
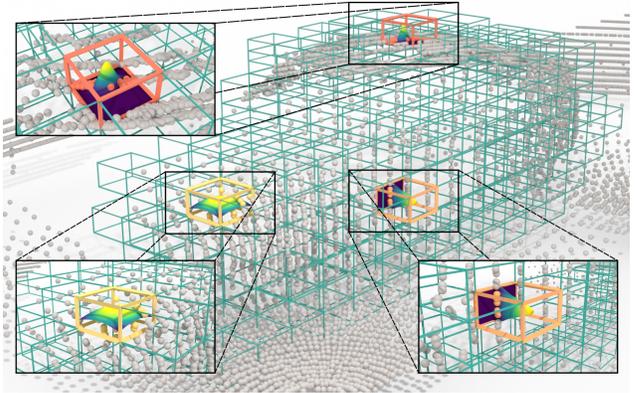


Figure 2. Illustration of our GBlobs encoding: We explicitly model the local geometry within point neighborhoods by Gaussian blobs (highlighted). Here, shown for a car's point cloud representation, all points within a voxel belong to the same neighborhood.

To mitigate these drawbacks, we propose modeling small local neighborhoods of LiDAR point clouds as individual Gaussian blobs, as illustrated in Fig. 2. Our approach offers several advantages: it preserves local structure more effectively than simple relative distance, thanks to its covariance modeling. Additionally, it is permutation-invariant, meaning it does not require point cloud sorting, and it avoids the overhead of surface normal calculations, resulting in exceptional efficiency.

### 3.1. Preliminaries

Formally, given a set of $D \geqslant 2$ domains $\{(X_d, Y_d)\}_{d=2}^D$, each with distinct data and label distributions, we train the model on a subset of $K$ domains and conduct the evaluation on a single domain. We adopt the prevalent assumption that the domain gap primarily stems from discrepancies in data distributions $X_d$, while label distributions $Y_d$ across domains are considered largely consistent. We then distinguish two different tasks and evaluation protocols: (1) Single Domain Generalization (SDG) [21, 28] when $K = 1$, and (2) Multi-Domain Generalization (MDG) [49] when $K = D - 1$.

A LiDAR point cloud is represented as an unordered set of points, typically defined as $\{p_j = (x, y, z, i)\}_{j=1}^M \sim X_d$, where each point $p_j$ has a 3D coordinate $(x, y, z)$ relative to the sensor's origin and an associated intensity value $(i)$. Prior research on DG consistently overlooked additional explicit and implicit LiDAR channels, relying solely on $(x, y, z)$ as model input. For instance, due to the inherent challenges posed by intensity data, it is commonly excluded from the input features [18]. Similarly, the influence of different input descriptors, like relative distances [8, 55], surface normals [37, 69], or local encoding [20, 55, 70], remains severely underexplored.

## 3.2. Local Point Clouds Structure as Gaussian Blobs

Previous approaches conveniently train a 3D detection model on $K$ domains $\{(X_d, Y_d)\}_{d=2}^{K}$, using global point positions, $\{p_j = (x, y, z)\}_{j=1}^{M} \sim X_d$, as input features. The models are subsequently evaluated on the remaining domain using the same input representation.

Training with such global input features biases the detector towards the source domain(s): Different reference coordinate systems (*e.g.* compare KITTI [12] *vs.* Waymo [50]) lead to different distributions of point coordinates, which are inevitably picked up during training and often hinder the generalization performance. Thus, it is common practice in both, DG [21, 49, 59] and UDA [5, 63] approaches, to shift the datasets in the height axis in order to manually align their coordinate distributions.

Conversely, we decouple the absolute object location from model encoding while preserving spatial information by representing a small local point cloud neighborhood as a Gaussian blob, $\mathcal{N}(\boldsymbol{\mu}, \Sigma)$, where

$$\boldsymbol{\mu} = \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{p}_i, \quad \text{and} \tag{1}$$

$$\Sigma = \frac{1}{N} \sum_{i=1}^{N} (\boldsymbol{p}_i - \boldsymbol{\mu})(\boldsymbol{p}_i - \boldsymbol{\mu})^{\top}. \tag{2}$$

The neighborhood of $N$ points is determined by the detector. For instance, voxel-based detectors naturally limit this neighborhood to the maximum points per voxel. The dimension of $\boldsymbol{\mu} \in \mathbb{R}^M$ and $\Sigma \in \mathbb{R}^{M \times M}$ depends on the dimension of the LiDAR points $p_j \in \mathbb{R}^M$. Although our representation can leverage intensity information by setting $M = 4$ during training, we conduct all our experiments using $M = 3$ to ensure a fair comparison with state-of-the-art DG methods. These methods typically exclude intensity information due to its adversarial impact on cross-domain performance, as noted in prior works [18, 49]. In total, the dimensionality of our feature vector is $M + M^2 = 12$.

While $\Sigma$ is decoupled from the absolute point position, $\boldsymbol{\mu}$ remains expressed in the LiDAR's coordinate reference frame. Consequently we simply express $\boldsymbol{\mu}$ in a local frame of reference

$$\boldsymbol{d} = \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{p}_i - \boldsymbol{\mu}. \tag{3}$$

We directly use $(\boldsymbol{d}, \Sigma)$ as an input to an arbitrary model. Algorithm 1 summarizes our Gaussian blob computation in pseudocode.

### 3.3. Global *vs.* Local Representations

The common practice of using global point cloud coordinates as input features in 3D models is analogous to appending pixel coordinates to RGB values of images. This, however,

**Algorithm 1** GBlobs in a PyTorch-like pseudocode.

```
# f: N queries, each containing K neighborhood points
    of dimension M (NxKxM)

# compute mean                              ▷ Eq. (1)
f_mean = f.mean(dim=1, keepdims=True)
# compute cov                               ▷ Eq. (2)
f_loc = f - f_mean
cov = torch.einsum("nka,nkb->nab", f_loc, f_loc)
# final descriptor
d = f_loc.mean(dim=1) # NxM                 ▷ Eq. (3)
cov_flat = cov.reshape(-1, pow(M, 2)) # NxM^2
gblobs = torch.cat([d, cov_flat], dim=1) # (Nx(M+M^2))
```

would violate the assumed translation invariance, because the same object at different locations would have distinct feature activations. Moreover, the explicit global representation is unnecessary, as positional information is already implicitly encoded in 2D/3D CNNs [31, 61]. Even transformer architectures [19, 30, 52, 55], which lack an inherent notion of position, explicitly inject it through different positional encodings. Therefore, the use of absolute position as additional input features is redundant and, as we show, can be detrimental for generalization performance.

Our local geometry $(\boldsymbol{d}, \Sigma)$ effectively separates an object's spatial position within a LiDAR frame from its encoding. Unlike global coordinates, similar local geometry induces similar activations, regardless of their absolute location. By explicitly representing the local geometry, we leave the positional information to the network architecture. This disentanglement allows the model to learn more abstract and invariant representations of objects, effectively replacing the redundant position information with a geometrically sound representation.

## 4. Experiments

We present detailed experiments to validate our approach on widely-used benchmarks. Sec. 4.1 and Sec. 4.2 present our findings from single-source [28] and multi-source [49] domain generalization experiments, respectively. Finally, Sec. 4.3 presents our thorough ablation studies.

**Datasets.** In order to be comparable with the existing benchmarks [21, 28, 49], we employ four commonly used autonomous driving datasets, as outlined in Tab. 1. The variety of datasets allows us to thoroughly validate our claims. Differences in the number of LiDAR beams and vertical field-of-view induce vastly different numbers of points per-sample (*i.e.* per LiDAR point cloud). Diverse LiDAR configurations create unique sampling patterns, biasing detection models towards their source data. This is evident in their poor cross-domain performance, as demonstrated in our subsequent evaluations.

**Metrics.** Following the established evaluation protocols [62, 63], we use the KITTI [12] evaluation metrics in all our experiments. If not stated otherwise, we utilize three main

| Dataset | Location | # Beams | PPS | VFOV |
|---------|----------|---------|-----|------|
| KITTI [12] | Germany | 64 | 118k | $[-23.6°, \ 3.2°]$ |
| Waymo [50] | USA | 64 | 160k | $[-17.6°, \ 2.4°]$ |
| nuScenes [2] | USA/Singapore | 32 | 25k | $[-30.0°, 10.0°]$ |
| ONCE [35] | China | 40 | 70k | $[-25.0°, 15.0°]$ |

Table 1. Characteristics of datasets used in our experiments, including their geographic location, number of LiDAR beams, points per-scan (PPS), and vertical field-of-view.

classes: Car (*i.e.* Vehicle for Waymo), Pedestrian and Cyclist. We report Average Precision (AP) computed over 40 recall points and mean AP (mAP) over all classes. The Intersection over Union (IoU) thresholds are set at 0.7, 0.5 and 0.5 for Car, Pedestrian and Cyclists, respectively.

**Baselines.** As a baseline, we report the performance of a model trained with default configuration, *i.e.*, its standard training settings (learning rate, number of epochs, *etc.*) and standard 3D data augmentations (ground truth sampling, rotation, scaling, and flipping). We compare this model with state-of-the-art single-source [7, 21, 28] and multi-source [49] domain generalization (DG) methods. To demonstrate our contributions, we train the default model again but replace the global inputs with our proposed method, as described in Sec. 3. For this, we do not employ any specialized data augmentation techniques, model parameter search, or hyperparameter optimization. For the evaluation, we simply select the last checkpoint.

### 4.1. Single-source Domain Generalization

Single-source domain generalization aims to train a model on a single source domain and generalize its performance to unseen target domains. In our experiments, we follow established benchmarks [21, 28, 62, 63] and test our method on three different dataset configurations: Waymo→KITTI, nuScenes→KITTI and Waymo→nuScenes. For apples-to-apples comparisons with other DG methods [7, 21, 28, 28], we use the same detector and point cloud configuration: We train Voxel R-CNN [8], limit the LiDAR range to $[-75.2m, -75.2m, -2m, 75.2m, 75.2m, 4m]$ and utilize a voxel size of $[0.1m, 0.1m, 0.15m]$. To train the model, we leveraged the OpenPCDet[1] framework.

We compare our method to PA-DA [7] and 3D-VF [21], which use different specialized data augmentation techniques to improve model generalization. Additionally, we also benchmark our method against the state-of-the-art by Li *et al.* [28], which uses a multi-task learning strategy in addition to a novel data augmentation technique. Moreover, they exploit a point cloud reconstruction task for test-time training. On the contrary, our method does not rely on heavy data augmentation, additional training tasks, or test-time training. We simply train the baseline model with

---

[1] https://github.com/open-mmlab/OpenPCDet

our GBlobs input features, using the default configuration. Our simple, yet highly effective approach surpasses existing DG approaches as shown in Tab. 2, achieving substantial improvements of over 21 and 12 mAP in Waymo→KITTI and nuScenes→KITTI, respectively. While the relatively small improvement observed in the Waymo→nuScenes experiment is discussed in detail in our ablation studies, the detector's inability to bridge the dense-to-sparse gap – a well-known challenge in 3D object detection [9, 58] – remains a primary limitation.

For additional insights, we replicate the KITTI→Waymo evaluation from 3D-VF [21]: We train three detectors (Point-Pillars [20], SECOND [61], and Part-A$^2$ [47]) using their default configurations (see supplementary material for details) and simply exchange their default global input features with our GBlobs. We evaluate these models on the full KITTI (in-domain) and Waymo eval split, using the Car/Vehicle class. Analogous to 3D-VF, we report AP with a 3D IoU of 0.7 for KITTI and 0.5 for Waymo. As shown in Tab. 3, our input encoding outperforms 3D-VF by over 4, 13, and 4 AP points using PointPillars, SECOND, and Part-A$^2$, respectively. Our proposed input encoding and standard detector settings were sufficient to achieve significant improvement without any specialized data augmentation. We observed that data augmentation tailored for the target domain can often negatively impact performance on the source domain, as demonstrated in the KITTI evaluation of Tab. 3. In contrast, our approach is unaffected by this issue and even outperforms other methods in in-domain evaluation.

### 4.2. Multi-source Domain Generalization

Multi-source DG trains models that can effectively generalize to unseen data from different distributions. By leveraging multiple source domains, the model can capture underlying patterns and invariant features that are common to all domains, enabling it to perform well on novel, unseen target domains. This approach is particularly valuable in scenarios where target domain data is limited or unavailable, making it a crucial tool for real-world applications, especially in fields like autonomous driving.

Given the limited research on LiDAR-based 3D object detection in multi-source domain generalization, we benchmark our approach against the state-of-the-art MDT3D [49]. Following their leave-one-out strategy (where one dataset is used for testing and the others for training), we train a CenterPoint [65] model on three of the four widely used autonomous driving datasets (KITTI [12], Waymo [50], nuScenes [2], and ONCE [2]) and evaluate on the remaining one. To ensure a fair comparison with MDT3D, we apply their uniform data subsampling schema and use the same LiDAR range $[-75.2m, -75.2m, -2m, 75.2m, 75.2m, 4m]$ and voxel size $[0.1m, 0.1m, 0.2m]$. Detailed experiment configurations are provided in the supplemental material.

| Tasks | Methods | Car | Pedestrian | Cyclist | mAP |
|---|---|---|---|---|---|
| Waymo → KITTI | Voxel R-CNN [8] | 66.65/19.27 | 66.55/64.00 | 63.04/57.11 | 65.41/46.79 |
| | PA-DA[7] | 65.82/17.61 | 66.40/63.88 | 61.30/56.23 | 64.51/45.91 |
| | 3D-VF[21] | 66.72/19.37 | 66.21/63.12 | 62.74/56.44 | 65.22/46.31 |
| | Li *et al.* [28] | 69.90/20.21 | 63.24/62.59 | 63.27/57.21 | 65.47/46.67 |
| | Voxel R-CNN w/ GBlobs | **87.33/78.75** | **69.47/65.98** | **63.63/58.72** | **73.48/67.82** |
| nuScenes → KITTI | Voxel R-CNN [8] | 66.93/28.80 | 23.39/18.65 | 19.23/15.76 | 36.52/21.07 |
| | PA-DA [7] | 65.09/32.44 | 18.73/14.94 | 18.66/15.91 | 34.16/21.10 |
| | 3D-VF [21] | 65.36/29.21 | 24.85/20.87 | 22.13/19.31 | 37.45/23.13 |
| | Li *et al.* [28] | 73.58/33.11 | 30.01/23.73 | 22.93/18.62 | 42.17/25.15 |
| | Voxel R-CNN w/ GBlobs | **80.95/53.98** | **38.33/33.22** | **29.18/25.68** | **49.48/37.62** |
| Waymo → nuScenes | Voxel R-CNN [8] | 31.20/19.13 | 10.52/ 8.39 | 0.75/ 0.55 | 14.16/ 9.36 |
| | PA-DA[7] | 29.43/18.06 | 10.84/ 8.43 | 0.82/ 0.43 | 13.70/ 8.97 |
| | 3D-VF[21] | 30.17/18.91 | 10.54/ 7.23 | 0.76/ 0.78 | 13.82/ 8.97 |
| | Li *et al.* [28] | **36.04/22.25** | **14.48/10.56** | 1.15/ 0.95 | **17.22**/11.26 |
| | Voxel R-CNN w/ GBlobs | 32.08/20.08 | 11.60/ 9.13 | **5.67/ 5.10** | 16.45/**11.44** |

Table 2. Single-source domain generalization experiments using Voxel R-CNN [8]. Following Li *et al.* [28], we trained a Voxel R-CNN detector on all three classes (Car/Vehicle, Pedestrian, Cyclist) simultaneously and evaluated performance using Average Precision (AP) on Bird's-eye View (BEV) / 3D views at 40 recall positions. Intersection over Union (IoU) thresholds of 0.7, 0.5, and 0.5 were used for Car/Vehicle, Pedestrian, and Cyclist, respectively. For KITTI evaluation, we report the average AP across all difficulty levels (Easy, Moderate, Hard). Additionally, we provide the mean AP over the three classes. The best value in each category is highlighted in bold.

| Model | Method | KITTI | Waymo |
|---|---|---|---|
| PointPil. [20] | Default | 77.11 | 40.86 |
| | 3D-VF [21] | 77.13 | 44.61 |
| | w/ GBlobs | **78.75** | **48.92** |
| SECOND [61] | Default | 78.68 | 42.45 |
| | 3D-VF [21] | 78.56 | 43.51 |
| | w/ GBlobs | **81.61** | **56.52** |
| Part-A$^2$ [47] | Default | 79.16 | 49.76 |
| | 3D-VF [21] | 79.26 | 56.08 |
| | w/ GBlobs | **82.29** | **60.20** |

Table 3. Following 3D-VF [21], we trained three detectors on the KITTI dataset (Car class only), and evaluated them on KITTI (Moderate difficulty, IoU threshold 0.7) and Waymo (IoU threshold 0.5) using 3D Average Precision (AP).

In Tab. 4, we report the mean Average Precision (mAP) computed over all classes for each dataset. Following [49], for KITTI, we report the moderate difficulty, whereas for the others there is no such categorization. Unlike MDT3D, which uses multi-domain data mixing, our method employs the standard CenterPoint configuration with default augmentations (ground truth sampling, rotation, scaling, and flipping) and no additional hyperparameter tuning. Without any bells and whistles our method exhibits tremendous improvements over MDT3D of around 17, 12 and 5 mAP for Waymo, KITTI and ONCE. Similar to Sec. 4.1, extremely

sparse neighborhoods (*i.e.* nuScenes) poses a challenge, as we further investigate within our ablation studies.

### 4.3. Ablation Studies

**Input Feature Impact.** In Tab. 5, we compare the performance of models trained on different input features for in- and cross-domain evaluation. Our results show that the choice of input features has little impact on the in-domain performance. However, local features significantly outperform global and hybrid (combined local and global) features in the cross-domain scenario. There, our proposed GBlobs, *i.e.* $(\mathbf{b}, \Sigma)$, achieve the best result.

**In-domain Performance.** While 3D-VF [21] found varying levels of in-domain performance among different 3D detectors (similar to our Tab. 3), Li *et al.* [28] and MDT3D [49] do not report any in-domain evaluation results. Given that their data augmentation strategies are focused on target datasets, it is likely that their model performance on the source dataset would degrade, although the extent is uncertain. While our Gaussian blobs outperform these methods in competitive domain generalization (DG) benchmarks, the impact of our approach on in-domain performance remains unclear. To assess this, we evaluated various detectors on different autonomous driving datasets, as shown in Table 6. Our experiments on the KITTI dataset (Tab. 6a) show that our input features surpass the default detector in all classes. Additionally, our method improves a DSVT (Pillar) transformer [55] by over 3 mAP on the Waymo dataset (Tab. 6b).

| Method | KITTI | ONCE | nuScenes | Waymo | mAP |
|---|---|---|---|---|---|
| KITTI | 47.40 | 9.90 | 1.60 | 3.00 | 4.83 |
| ONCE | 43.00 | 60.20 | 5.80 | 13.00 | 20.60 |
| nuScenes | 12.60 | 10.70 | 18.40 | 0.90 | 8.06 |
| Waymo | 34.20 | 28.80 | 7.20 | 40.20 | 23.40 |
| MDT3D [49] | 41.80 | 32.28$^{\dagger}$ | **11.00** | 6.40 | 22.87 |
| CP w/ GBlobs | **53.92** | **37.77** | 8.15 | **23.81** | **30.91** |

Table 4. Multi-source DG results. The top part presents the results for a vanilla CenterPoint [65] (CP) detector, trained exclusively on single domains (leftmost column) and evaluated on all datasets. For each dataset, we report mean Average Precision (mAP) over all classes (columns KITTI–Waymo) and across all datasets, excluding in-domain experiments (rightmost column). The bottom part presents the multi-source DG results following MDT3D's leave-one-out strategy, *i.e.* CenterPoint is trained with our GBlobs as input features on three datasets and evaluated on the fourth. $^{\dagger}$: in Tab. 7 of MDT3D [49], the mAP is wrongly calculated; we report the corrected value.

| global | d | Σ | n→n | n→K |
|---|---|---|---|---|
| ✓ | | | **30.21/22.49** | 36.52/21.07 |
| ✓ | | ✓ | 29.79/22.37 | 36.58/22.45 |
| | ✓ | | 29.29/21.84 | 48.61/35.95 |
| | ✓ | ✓ | 28.51/21.13 | 49.02/37.60 |
| ✓ | ✓ | ✓ | 29.33/21.90 | **49.48/37.62** |

Table 5. Impact of different input features on Voxel R-CNN [8]. We evaluate the in-domain nuScenes→nuScenes (n→n) and cross-domain nuScenes→KITTI (n→K) performance. We report mAP, computed over all three classes, on BEV / 3D views at 40 recall positions. We denote standard training with global features (*i.e.* Cartesian coordinates) as "global". d and Σ are features computed as in Eq. (3) and Eq. (2), respectively.

| Method | Car | Pedestrian | Cyclist | mAP |
|---|---|---|---|---|
| SEC. [61] | 88.2/80.9 | 54.2/49.8 | 68.0/63.1 | 70.1/64.6 |
| SEC. w/ GBlobs | **88.8/81.0** | **55.7/50.4** | **69.0/64.5** | **71.2/65.3** |

(a) In-domain evaluation of SECOND [61] on the KITTI *val* split (Moderate case), reporting Average Precision (AP) BEV/3D for IoU thresholds of 0.7 (Car) and 0.5 (Pedestrian and Cyclist).

| Method | Vehicle | Pedestrian | Cyclist | mAP |
|---|---|---|---|---|
| DSVT [55] | 74.1/65.7 | 73.3/65.4 | 58.0/55.8 | 68.5/62.3 |
| DSVT w/ GBlobs | **75.4/67.0** | **76.7/68.8** | **64.4/62.0** | **72.2/65.9** |

(b) In-domain evaluation on the full Waymo *val* split reporting L1/L2 Average Precision (AP). Models are trained on 20% of the training split.

Table 6. In-domain experiments on KITTI [12] and Waymo [50].

**Different Local Descriptors.** A limited subset of methods augment point clouds with local geometric features, such as relative distance from voxels [8] or umbrella-like hand-crafted descriptors [43], in addition to standard global inputs. However, models trained on such hybrid features inherit all the disadvantages of global geometric features (see Tab. 5). Therefore, in Tab. 7, we remove all global components (if there are any) and compare various local point cloud features in the Waymo→KITTI case.

| Method | Car | Pedestrian | Cyclist | mAP |
|---|---|---|---|---|
| rel. distance [8] | 84.93/74.60 | 68.72/65.46 | 60.47/55.96 | 71.38/65.34 |
| surf. normals [15] | 85.39/73.37 | 66.32/63.96 | **63.96**/58.44 | 70.05/63.96 |
| RepSurf [43] | 84.45/73.47 | 69.41/**66.87** | 61.98/57.70 | 71.95/66.02 |
| GBlobs | **87.33/78.75** | **69.47**/65.98 | 63.63/**58.72** | **73.48/67.82** |

Table 7. Waymo→KITTI evaluation of different local descriptors using Voxel R-CNN [8]. We report AP on BEV / 3D views at 40 recall positions and mAP computed over all three classes.

Encoding local geometry simply as relative distances is invariant to global transformations and provides good results. However, this encoding is limited in its ability to capture fine-grained details, which handcrafted descriptors (*e.g.* RepSurf [43]) can provide. Unlike RepSurf, our method does not require surface normal estimates, which are highly sensitive to noise [45]. Moreover, it is inherently permutational invariant and does not require any point cloud sorting, yet it consistently outperforms all other methods on the Waymo→KITTI DG benchmark.

**Information Retention.** Voxel-based detectors (*e.g.* [8, 61]) average points within voxels to create per-voxel features, which can lead to data loss. Using smaller voxels mitigates this by increasing granularity but at a higher computational cost. Pillar-based detectors (*e.g.* [20, 24]) use voxels spanning the full height, often employing a small PointNet [40] to encode pillar features, though some information loss is inevitable. Our Gaussian blobs outperform existing encoders in preserving information: as shown in Fig. 3, larger voxels degrade performance for both in-domain (Fig. 3a) and cross-domain (Fig. 3b) evaluations, regardless of the encoder. Our encoding consistently surpasses PointNet voxel encoders, even with fewer parameters, likely because voxel encoders struggle to reliably encode local structure from global coordinates without additional supervision.

**nuScenes Performance.** Both single-source (Tab. 2) and multi-source (Tab. 4) DG approaches struggle with the challenging nuScenes benchmark. This dense-to-sparse gap is a
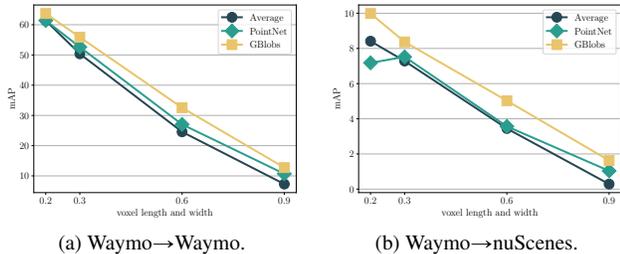
(a) Waymo→Waymo.  (b) Waymo→nuScenes.

Figure 3. Impact of voxel size and voxel feature encoder for Voxel R-CNN [8], trained on Waymo and evaluated for the (a) in-domain and (b) cross-domain setting. Best viewed on screen.



(a) In-domain: SECOND [61] on KITTI→"sparser KITTI".  (b) Cross-domain: Voxel R-CNN [8] on nuScenes→"sparser KITTI".
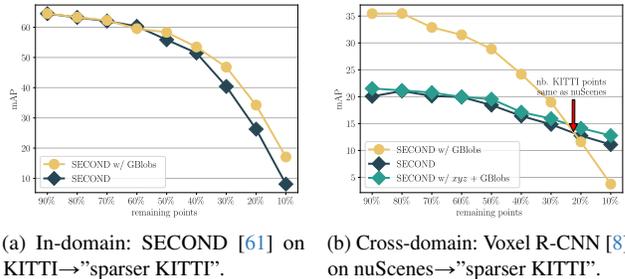
Figure 4. Extreme sparsity evaluation: detectors trained on complete point clouds were evaluated on gradually sub-sampled KITTI point clouds at test time, both for the (a) in-domain and (b) cross-domain setting. Best viewed on screen.

well-known challenge in 3D object detection [9, 58]. While our GBlobs consistently yields significant improvements across all other benchmarks, the achievable performance on nuScenes is on-par with default global input features.

Compared to Waymo and KITTI, nuScenes has significantly fewer points per LiDAR frame (recall Tab. 1) and a distribution that is skewed towards the lower range: only about 1.5k voxels (6% of the total) contain more than three points, while 94% of all nuScenes voxels contain only one or two points. The underdetermined system, lacking sufficient data to estimate the covariance matrix as per Eq. (2), results in a degenerate mean-only case. This is because the covariance matrix values are forced to zero, significantly limiting the model's ability to capture complex data distributions and leading to poor nuScenes performance. However, we anticipate that with the constant sensor improvements, such sparse LiDAR point-clouds (e.g. sparser than nuScenes) will soon become obsolete, making this generalization scenario highly unlikely in real-world autonomous driving applications.

**Sparsity Influence.** To evaluate the impact of sparsity on our input representation, we conduct two experiments: an in-domain evaluation on KITTI [12] using SECOND [61] and a cross-domain nuScenes→KITTI evaluation using Voxel R-CNN [8]. The in-domain model was trained using the standard KITTI configuration, while the cross-domain experiment followed the setup described in Sec. 4.1. In both cases, we report mean Average Precision (mAP) computed over the Car, Pedestrian, and Cyclist classes.

During testing, we uniformly sampled the input point cloud to retain a specified portion of the points, discarding the rest. As illustrated in Fig. 4, a higher point cloud density correlates with improved precision for both in-domain and cross-domain experiments. Conversely, reducing the number of input points resulted in a decrease in performance.

In Fig. 4a, reducing the number of input points transforms the task into a domain generalization problem. In such cases, as previously shown, our method excels. Even when 90% of the data points were discarded, the model trained with our inputs outperformed the original by over 10 mAP.

Our sparse cross-domain experiment in Fig. 4b reveals

that in the extreme "sparse-to-sparser" case, where the target dataset has far fewer points than the already sparse source dataset, global point cloud features outperform local features. At around 25% remaining points, a subsampled KITTI frame has roughly the same number of points as nuScenes, indicating that global location data should be preferred when a detector is pre-trained on sparse data and local neighborhoods are unreliable. However, this scenario is unlikely in real-world applications. Notably, as sparsity increases, the combined model ($xyz$ + GBlobs) degrades less significantly than GBlobs alone, suggesting it leverages GBlobs in dense voxels while relying on global coordinates elsewhere.

## 5. Conclusion

In this work, we emphasized the critical role of local point cloud geometry in building robust LiDAR-based 3D detection models. Our evaluations revealed that relying solely on standard global input features leads to poor cross-domain generalization. To address this limitation, we introduced a highly efficient, parameter-free approach that can be seamlessly integrated into any 3D object detector: encoding local point cloud geometry using Gaussian blobs. Our method, GBlobs, has been extensively validated across a range of detectors and benchmark scenarios, demonstrating competitive in-domain performance (on par with standard global features) while achieving significant improvements in cross-domain generalization.

## References

[1] Sravanti Addepalli, Ashish Ramayee Asokan, Lakshay Sharma, and R Venkatesh Babu. Leveraging Vision-Language

Models for Improving Domain Generalization in Image Classification. In *Proc. CVPR*, 2024. 3

[2] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuScenes: A multimodal dataset for autonomous driving. In *Proc. CVPR*, 2020. 5, 12, 13

[3] Benjamin Caine, Rebecca Roelofs, Vijay Vasudevan, Jiquan Ngiam, Yuning Chai, Zhifeng Chen, and Jonathon Shlens. Pseudo-labeling for Scalable 3D Object Detection. *arXiv CoRR*, abs/2103.02093, 2021. 2

[4] Yukang Chen, Jianhui Liu, Xiangyu Zhang, Xiaojuan Qi, and Jiaya Jia. VoxelNeXt: Fully Sparse VoxelNet for 3D Object Detection and Tracking. In *Proc. CVPR*, 2023. 1, 2

[5] Zhuoxiao Chen, Yadan Luo, Zheng Wang, Mahsa Baktashmotlagh, and Zi Huang. Revisiting Domain-Adaptive 3D Object Detection by Reliable, Diverse and Class-balanced Pseudo-Labeling. In *Proc. ICCV*, 2023. 2, 4

[6] Zining Chen, Weiqiu Wang, Zhicheng Zhao, Fei Su, Aidong Men, and Hongying Meng. PracticalDG: Perturbation Distillation on Vision-Language Models for Hybrid Domain Generalization. In *Proc. CVPR*, 2024. 3

[7] Jaeseok Choi, Yeji Song, and Nojun Kwak. Part-Aware Data Augmentation for 3D Object Detection in Point Cloud. In *Proc. IROS*, 2021. 5, 6, 13

[8] Jiajun Deng, Shaoshuai Shi, Peiwei Li, Wengang Zhou, Yanyong Zhang, and Houqiang Li. Voxel R-CNN: Towards High Performance Voxel-based 3D Object Detection. In *Proc. AAAI*, 2021. 1, 2, 3, 5, 6, 7, 8, 13, 14

[9] George Eskandar. An Empirical Study of the Generalization Ability of Lidar 3D Object Detectors to Unseen Domains. In *Proc. CVPR*, 2024. 1, 5, 8

[10] Lue Fan, Ziqi Pang, Tianyuan Zhang, Yu-Xiong Wang, Hang Zhao, Feng Wang, Naiyan Wang, and Zhaoxiang Zhang. Embracing Single Stride 3D Object Detector with Sparse Transformer. In *Proc. CVPR*, 2022. 2

[11] Christian Fruhwirth-Reisinger, Michael Opitz, Horst Possegger, and Horst Bischof. FAST3D: Flow-Aware Self-Training for 3D Object Detectors. In *Proc. BMVC*, 2021. 2

[12] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite. In *Proc. CVPR*, 2012. 4, 5, 7, 8, 12, 13

[13] Rui Gong, Wen Li, Yuhua Chen, and Luc Van Gool. DLOW: Domain Flow for Adaptation and Generalization. In *Proc. CVPR*, 2019. 3

[14] Martin Hahner, Dengxin Dai, Alexander Liniger, and Luc Van Gool. Quantifying Data Augmentation for LiDAR-based 3D Object Detection. *arXiv CoRR*, abs/2004.01643, 2020. 3

[15] Hugues Hoppe, Tony DeRose, Tom Duchamp, John McDonald, and Werner Stuetzle. Surface Reconstruction from Unorganized Points. In *SIGGRAPH*, 1992. 2, 3, 7

[16] Qianjiang Hu, Daizong Liu, and Wei Hu. Density-Insensitive Unsupervised Domain Adaption on 3D Object Detection. In *Proc. CVPR*, 2023. 2

[17] Jiaxing Huang, Dayan Guan, Aoran Xiao, and Shijian Lu. FSDR: Frequency Space Domain Randomization for Domain Generalization. In *Proc. CVPR*, 2021. 3

[18] Jaeyeul Kim, Jeonghoon Woo, Jungwan sand Kim, and Sunghoon Im. Rethinking LiDAR Domain Generalization: Single Source as Multiple Density Domains. In *Proc. ECCV*, 2024. 3, 4

[19] Alexander Kolesnikov, Alexey Dosovitskiy, Dirk Weissenborn, Georg Heigold, Jakob Uszkoreit, Lucas Beyer, Matthias Minderer, Mostafa Dehghani, Neil Houlsby, Sylvain Gelly, Thomas Unterthiner, and Xiaohua Zhai. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *Proc. ICML*, 2021. 4

[20] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. PointPillars: Fast Encoders for Object Detection from Point Clouds. In *Proc. CVPR*, 2019. 2, 3, 5, 6, 7, 13

[21] Alexander Lehner, Stefano Gasperini, Alvaro Marcos-Ramiro, Michael Schmidt, Mohammad-Ali Nikouei Mahani, Nassir Navab, Benjamin Busam, and Federico Tombari. 3D-VField: Adversarial Augmentation of Point Clouds for Domain Generalization in 3D Object Detection. In *Proc. CVPR*, 2022. 2, 3, 4, 5, 6, 12

[22] Zhaoqi Leng, Guowang Li, Chenxi Liu, Ekin Dogus Cubuk, Pei Sun, Tong He, Dragomir Anguelov, and Mingxing Tan. LidarAugment: Searching for Scalable 3D LiDAR Data Augmentations. In *Proc. ICRA*, 2023. 3

[23] Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C Kot. Domain Generalization with Adversarial Feature Learning. In *Proc. CVPR*, 2018. 3

[24] Jinyu Li, Chenxu Luo, and Xiaodong Yang. PillarNeXt: Rethinking network designs for 3D object detection in LiDAR point clouds. In *Proc. CVPR*, 2023. 2, 7

[25] Li Li, Tanqiu Qiao, Hubert PH Shum, and Toby P Breckon. TraIL-Det: Transformation-Invariant Local Feature Networks for 3D LiDAR Object Detection with Unsupervised Pre-Training. In *Proc. BMVC*, 2024. 1, 2

[26] Li Li, Hubert PH Shum, and Toby P Breckon. RAPiD-Seg: Range-Aware Pointwise Distance Distribution Networks for 3D LiDAR Segmentation. In *Proc. ECCV*, 2024. 2

[27] Qing Li, Huifang Feng, Kanle Shi, Yue Gao, Yi Fang, Yu-Shen Liu, and Zhizhong Han. NeuralGF: Unsupervised Point Normal Estimation by Learning Neural Gradient Function. In *Proc. NeurIPS*, 2023. 3

[28] Shuangzhi Li, Lei Ma, and Xingyu Li. Domain Generalization of 3D Object Detection by Density-Resampling. In *Proc. ECCV*, 2024. 2, 3, 4, 5, 6

[29] Ya Li, Mingming Gong, Xinmei Tian, Tongliang Liu, and Dacheng Tao. Domain generalization via conditional invariant representations. In *Proc. AAAI*, 2018. 3

[30] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin Transformer: Hierarchical Vision Transformer using Shifted Windows. In *Proc. ICCV*, 2021. 4

[31] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A ConvNet for the 2020s. In *Proc. CVPR*, 2022. 4

[32] Fangrui Lv, Jian Liang, Shuang Li, Bin Zang, Chi Harold Liu, Ziteng Wang, and Di Liu. Causality Inspired Representation Learning for Domain Generalization. In *Proc. CVPR*, 2022. 3

[33] Divyat Mahajan, Shruti Tople, and Amit Sharma. Domain Generalization using Causal Matching. In *Proc. ICML*, 2021. 3

[34] Dušan Malić, Christian Fruhwirth-Reisinger, Horst Possegger, and Horst Bischof. SAILOR: Scaling Anchors via Insights into Latent Object Representation. In *Proc. WACV*, 2023. 2

[35] Jiageng Mao, Minzhe Niu, Chenhan Jiang, Hanxue Liang, Jingheng Chen, Xiaodan Liang, Yamin Li, Chaoqiang Ye, Wei Zhang, Zhenguo Li, et al. One Million Scenes for Autonomous Driving: ONCE Dataset. *Proc. NeurIPS*, 2022. 5

[36] Gregory P Meyer, Ankit Laddha, Eric Kee, Carlos Vallespi-Gonzalez, and Carl K Wellington. LaserNet: An Efficient Probabilistic 3D Object Detector for Autonomous Driving. In *Proc. CVPR*, 2019. 2

[37] Jishu Miao, Tsubasa Hirakawa, Takayoshi Yamashita, and Hironobu Fujiyoshi. 3D Object Detection with Normal-map on Point Clouds. In *Proc. VISIGRAPP*, 2021. 1, 2, 3

[38] Ammar Yasir Naich and Jesús Requena Carrión. LiDAR-Based Intensity-Aware Outdoor 3D Object Detection. *Sensors*, 24(9):2942, 2024. 1

[39] Xingchao Peng, Zijun Huang, Ximeng Sun, and Kate Saenko. Domain Agnostic Learning with Disentangled Representations. In *Proc. ICML*, 2019. 3

[40] Charles R. Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation. In *Proc. CVPR*, 2017. 1, 2, 7

[41] Charles R. Qi, Wei Liu, Chenxia Wu, Hao Su, and Leonidas J Guibas. Frustum PointNets for 3D Object Detection from RGB-D Data. In *Proc. CVPR*, 2018. 2

[42] Alexandre Rame, Corentin Dancette, and Matthieu Cord. Fishr: Invariant Gradient Variances for Out-of-Distribution Generalization. In *Proc. ICML*, 2022. 3

[43] Haoxi Ran, Jun Liu, and Chengjie Wang. Surface Representation for Point Clouds. In *Proc. CVPR*, 2022. 2, 3, 7

[44] Christoph B Rist, Markus Enzweiler, and Dariu M Gavrila. Cross-Sensor Deep Domain Adaptation for LiDAR Detection and Segmentation. In *Proc. IV*, 2019. 2

[45] Dominik Scheuble, Chenyang Lei, Seung-Hwan Baek, Mario Bijelic, and Felix Heide. Polarization Wavefront Lidar: Learning Large Scene Reconstruction from Polarized Wavefronts. In *Proc. CVPR*, 2024. 7

[46] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. PointR-CNN: 3D Object Proposal Generation and Detection from Point Cloud. In *Proc. CVPR*, 2019. 1, 2

[47] Shaoshuai Shi, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. From Points to Parts: 3D Object Detection From Point Cloud With Part-Aware and Part-Aggregation Network. *TPAMI*, 43(8):2647–2664, 2020. 5, 6

[48] Shaoshuai Shi, Li Jiang, Jiajun Deng, Zhe Wang, Chaoxu Guo, Jianping Shi, Xiaogang Wang, and Hongsheng Li. PV-RCNN++: Point-Voxel Feature Set Abstraction With Local Vector Representation for 3D Object Detection. *IJCV*, 131 (2):531–551, 2023. 1

[49] Louis Soum-Fontez, Jean-Emmanuel Deschaud, and François Goulette. MDT3D: Multi-Dataset Training for LiDAR 3D Object Detection Generalization. In *Proc. IROS*, 2023. 2, 3, 4, 5, 6, 7

[50] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in Perception for Autonomous Driving: Waymo Open Dataset. In *Proc. CVPR*, 2020. 4, 5, 7, 12, 13

[51] Pei Sun, Mingxing Tan, Weiyue Wang, Chenxi Liu, Fei Xia, Zhaoqi Leng, and Dragomir Anguelov. SWFormer: Sparse Window Transformer for 3D Object Detection in Point Clouds. In *Proc. ECCV*, 2022. 2

[52] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. In *Proc. NeurIPS*, 2017. 4

[53] Riccardo Volpi and Vittorio Murino. Addressing Model Vulnerability to Distributional Shifts over Image Transformation Sets. In *Proc. ICCV*, 2019. 3

[54] Sean Walsh, Jason Ku, Alex D. Pon, and Steven L. Waslander. Leveraging Temporal Data for Automatic Labelling of Static Vehicles. In *Proc. CVPR*, 2020. 2

[55] Haiyang Wang, Chen Shi, Shaoshuai Shi, Meng Lei, Sen Wang, Di He, Bernt Schiele, and Liwei Wang. DSVT: Dynamic Sparse Voxel Transformer with Rotated Sets. In *Proc. CVPR*, 2023. 1, 2, 3, 4, 6, 7, 13

[56] Yan Wang, Xiangyu Chen, Yurong You, Li Erran Li, Bharath Hariharan, Mark Campbell, Kilian Q Weinberger, and Wei-Lun Chao. Train in Germany, Test in The USA: Making 3D Object Detectors Generalize. In *Proc. CVPR*, 2020. 1, 2

[57] Zhaoxuan Wang, Xu Han, Hongxin Liu, and Xianzhi Li. RIDE: Boosting 3D Object Detection for LiDAR Point Clouds via Rotation-Invariant Analysis. *arXiv CoRR*, abs/2408.15643, 2024. 1

[58] Yi Wei, Zibu Wei, Yongming Rao, Jiaxin Li, Jie Zhou, and Jiwen Lu. LiDAR Distillation: Bridging the Beam-Induced Domain Gap for 3D Object Detection. In *Proc. ECCV*, 2022. 2, 5, 8

[59] Guile Wu, Tongtong Cao, Bingbing Liu, Xingxin Chen, and Yuan Ren. Towards Universal LiDAR-Based 3D Object Detection by Multi-Domain Knowledge Transfer. In *Proc. ICCV*, 2023. 2, 3, 4

[60] Xiaoyang Wu, Li Jiang, Peng-Shuai Wang, Zhijian Liu, Xihui Liu, Yu Qiao, Wanli Ouyang, Tong He, and Hengshuang Zhao. Point Transformer V3: Simpler, Faster, Stronger. In *Proc. CVPR*, 2024. 2

[61] Yan Yan, Yuxing Mao, and Bo Li. SECOND: Sparsely Embedded Convolutional Detection. *Sensors*, 18(10):3337, 2018. 1, 2, 4, 5, 6, 7, 8, 12, 13, 15

[62] Jihan Yang, Shaoshuai Shi, Zhe Wang, Hongsheng Li, and Xiaojuan Qi. ST3D: Self-training for Unsupervised Domain Adaptation on 3D Object Detection. In *Proc. CVPR*, 2021. 1, 2, 4, 5, 12

[63] Jihan Yang, Shaoshuai Shi, Zhe Wang, Hongsheng Li, and Xiaojuan Qi. ST3D++: Denoised Self-Training for Unsupervised Domain Adaptation on 3D Object Detection. *TPAMI*, 45(5):6354–6371, 2022. 1, 2, 4, 5, 12

[64] Zetong Yang, Yanan Sun, Shu Liu, and Jiaya Jia. 3DSSD: Point-based 3D Single Stage Object Detector. In *Proc. CVPR*, 2020. 2

[65] Tianwei Yin, Xingyi Zhou, and Philipp Krahenbuhl. Center-based 3D Object Detection and Tracking. In *Proc. CVPR*, 2021. 1, 2, 5, 7, 13

[66] Yurong You, Carlos Andres Diaz-Ruiz, Yan Wang, Wei-Lun Chao, Bharath Hariharan, Mark Campbell, and Kilian Q Weinbergert. Exploiting Playbacks in Unsupervised Domain Adaptation for 3D Object Detection. In *Proc. ICRA*, 2022. 2

[67] Xiangyu Yue, Yang Zhang, Sicheng Zhao, Alberto Sangiovanni-Vincentelli, Kurt Keutzer, and Boqing Gong. Domain Randomization and Pyramid Consistency: Simulation-to-Real Generalization without Accessing Target Domain Data. In *Proc. ICCV*, 2019. 3

[68] Zhanwei Zhang, Minghao Chen, Shuai Xiao, Liang Peng, Hengjia Li, Binbin Lin, Ping Li, Wenxiao Wang, Boxi Wu, and Deng Cai. Pseudo Label Refinery for Unsupervised Domain Adaptation on Cross-dataset 3D Object Detection. In *Proc. CVPR*, 2024. 2

[69] Yan Zhao, Jihong Zhu, Haoyu Liang, and Lyujie Chen. Integration of Coordinate and Geometric Surface Normal for 3D Point Cloud Object Detection. In *Proc. IJCNN*, 2021. 1, 2, 3

[70] Yin Zhou and Oncel Tuzel. VoxelNet: End-to-End Learning for Point Cloud Based 3D Object Detection. In *Proc. CVPR*, 2018. 2, 3

# GBlobs: Explicit Local Structure via Gaussian Blobs for Improved Cross-Domain LiDAR-based 3D Object Detection

## Supplementary Material

This supplementary material details the experimental setup (App. A), presents further findings on the influence of local and global features on LiDAR-based 3D object detection (App. B), and provides a qualitative analysis of our results (App. C).

## A. Experimental Setup

To ensure reproducibility (besides the provided source code), we present detailed information about our experimental setup in Tab. 8. If not specified otherwise, we use default settings from OpenPCDet[2].

Our models were trained on the entire KITTI and nuScenes training sets, along with $20\%$ of the Waymo dataset (a standard practice in the field). In all our experiments (except KITTI→Waymo in Tab. 3 of the main manuscript), we trained the models to simultaneously predict Cars/Vehicles, Pedestrians, and Cyclists. For fair comparison with 3D-VF [21], we trained the detector to predict only Cars/Vehicles in KITTI→Waymo. We employed standard data augmentation techniques, including random sampling, point cloud rotation, scaling, and flipping.

We use the KITTI metric [12] for evaluation (except in Tab. 6b of the main manuscript), reporting Average Precision (AP) on Bird's-eye View (BEV) / 3D views at $40$ recall positions. For the in-domain Waymo evaluation in Tab. 6b (main manuscript), we report LEVEL_1 /LEVEL_2 AP (standard Waymo metric). We use Intersection over Union (IoU) thresholds of $0.7$, $0.5$, and $0.5$ for Cars, Pedestrians, and Cyclists, respectively. KITTI→Waymo in Tab. 3 (main manuscript) uses an IoU threshold of $0.5$ for Cars to ensure fair comparison. We utilize the complete validation sets of all datasets to assess the performance of our proposed method.

## B. Height Bias

Autonomous driving datasets define different reference points for LiDAR point clouds, *e.g.* Waymo [50] aligns the height axis origin with the road, while KITTI [12] and nuScenes [2] use the vehicle's mounting point. This inherently introduces bias into the network. A common approach is to manually align source and target point clouds by shifting them to a shared origin [62, 63]. Otherwise the detectors fail catastrophically as demonstrated in Tab. 9. Although this is not a critical issue in our controlled setting, a detector trained with such bias could pose a significant risk in real-world applications. Our GBlobs are not affected by biases
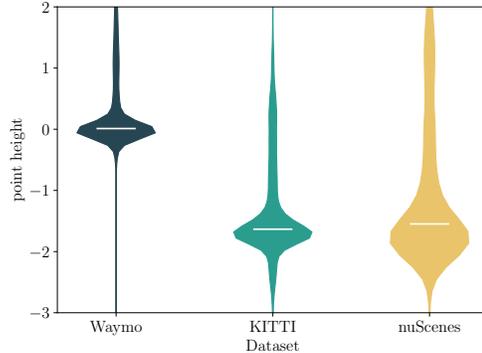


Figure 5. $z$ distribution

associated with global input features, as they encodes local point cloud geometry.

## C. Qualitative Results

In order to depict benefits of training a model with our GBlobs as input features, we conduct following qualitative analysis. We apply a nuScenes trained Voxel R-CNN detector to a challenging KITTI scene featuring a slightly curved road. Such detectors, trained with standard global input features, often predict false positives, even in areas without object indications.

A similar phenomenon can be observed with the SECOND [61] detector employed in the KITTI→Waymo benchmark in Fig. 7. It is noteworthy that SECOND, trained on KITTI, a dataset consisting primarily of small and mid-size European sedans, has never seen anything that resembles aerial work platforms during training. Nevertheless, when trained with global features and applied on Waymo (which has such object labeled as Vehicles), it manages to produce a detection at this location with high certainty (orange arrows in Fig. 7a). We hypothesize that the detector's prediction was influenced by specific points at specific heights. Given its training, such detections are unexpected. This raises the question of how many other detections, which are false positives, such detector produces. A model trained with our GBlobs did not make such uneducated guess (Fig. 7b).
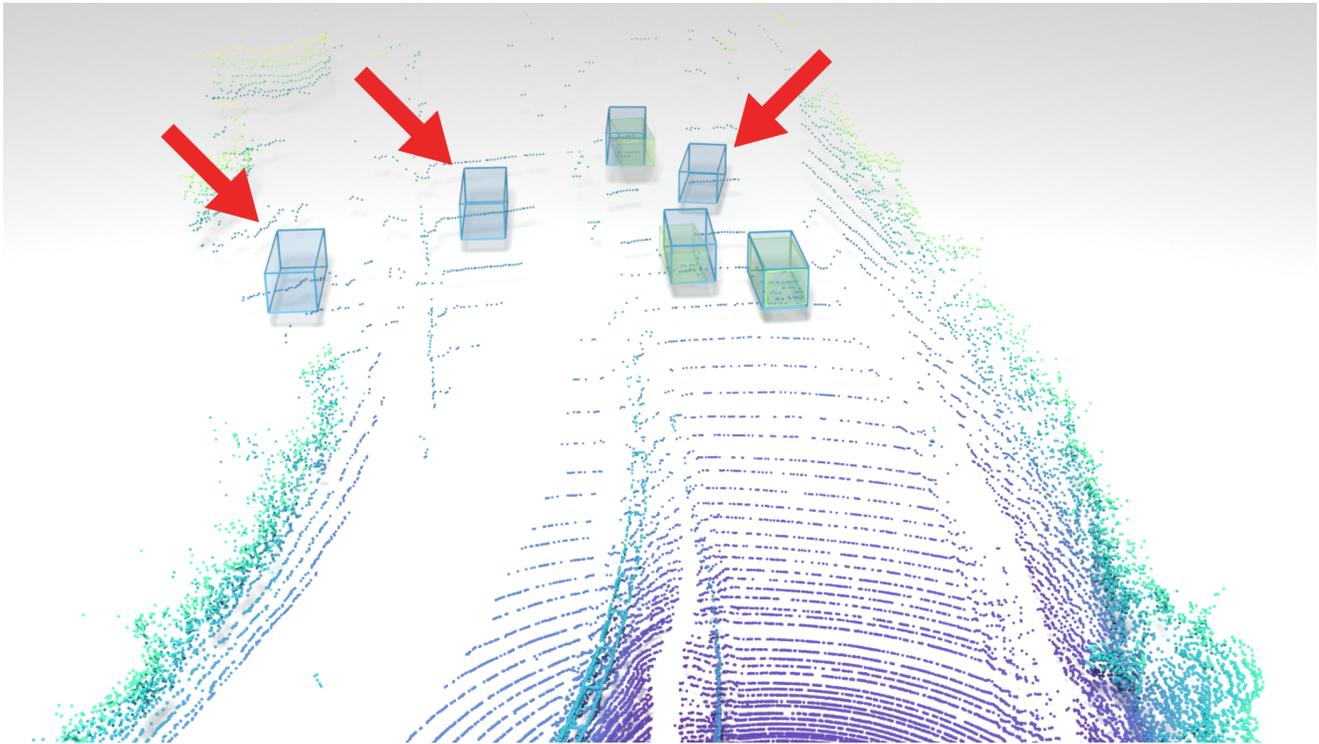
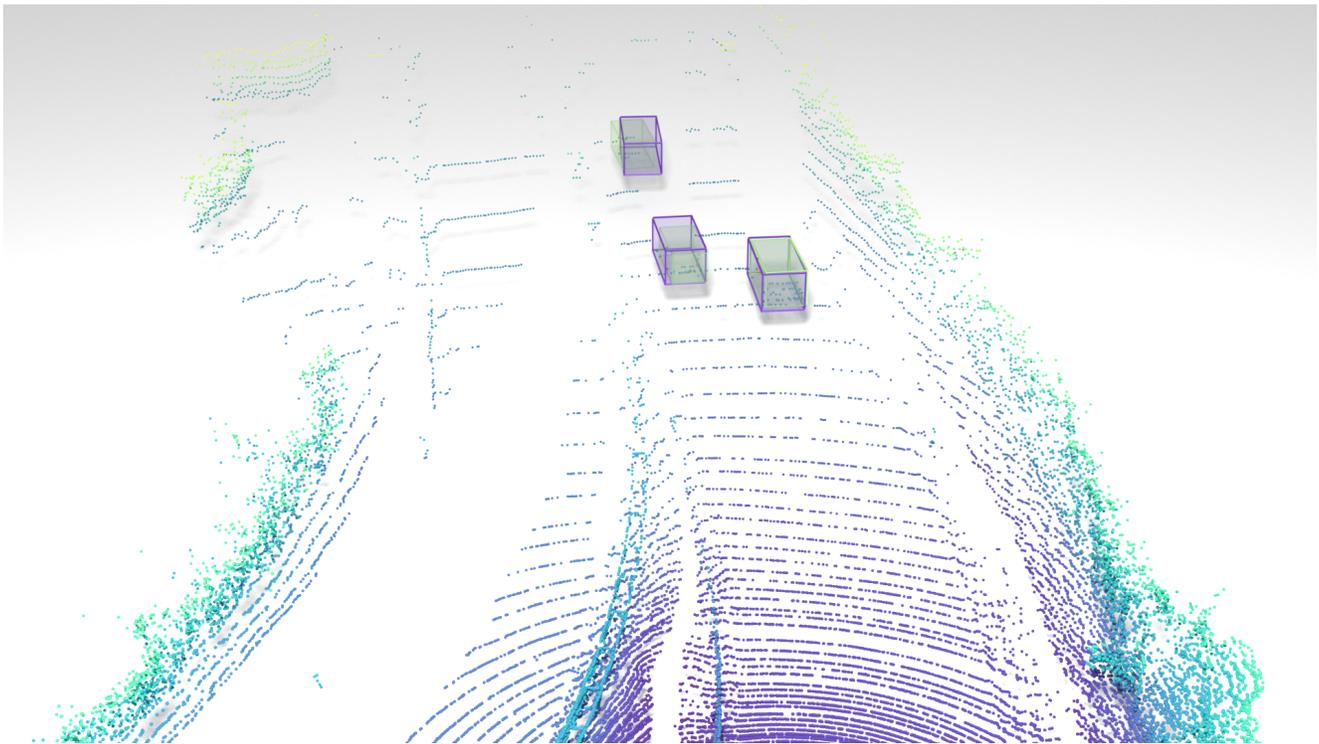| Table | Dataset | Detector | | | Optimizer | | | BS | E |
|---|---|---|---|---|---|---|---|---|---|
| | | Name | Range | Voxel Size | Name | LR | WD | | |
| Tab. 2 | Waymo [50] | Voxel R-CNN [8] | $[-75.2\ ,-75.2\ ,-2,75.2\ ,75.2\ ,4]$ | $[0.1\ ,0.1\ ,0.15]$ | Adam | $1\times10^{-2}$ | $1\times10^{-3}$ | 32 | 30 |
| | nuScenes [2] | Voxel R-CNN [8] | $[-75.2\ ,-75.2\ ,-2,75.2\ ,75.2\ ,4]$ | $[0.1\ ,0.1\ ,0.15]$ | Adam | $1\times10^{-2}$ | $1\times10^{-3}$ | 32 | 30 |
| Tab. 3 | KITTI [12] | PointPillars [20] | $[\ \ \ 0.0\ ,-39.68,-2,69.12,39.68,4]$ | $[0.16,0.16,6.0\ ]$ | Adam | $3\times10^{-3}$ | $1\times10^{-2}$ | 32 | 80 |
| | KITTI [12] | SECOND [61] | $[\ \ \ 0.0\ ,-40.0\ ,-3,70.4\ ,40.0\ ,1]$ | $[0.16,0.16,6.0\ ]$ | Adam | $3\times10^{-3}$ | $1\times10^{-2}$ | 32 | 80 |
| | KITTI [12] | Part-A$^2$ [7] | $[\ \ \ 0.0\ ,-40.0\ ,-3,70.4\ ,40.0\ ,1]$ | $[0.16,0.16,6.0\ ]$ | Adam | $1\times10^{-2}$ | $1\times10^{-2}$ | 32 | 80 |
| Tab. 4 | $*$ | CenterPoint [65] | $[-75.2\ ,-75.2\ ,-3,75.2\ ,75.2\ ,5]$ | $[0.10,0.10,0.20]$ | Adam | $3\times10^{-3}$ | $1\times10^{-2}$ | 32 | 30 |
| Tab. 5 | nuScenes [2] | Voxel R-CNN [8] | $[-75.2\ ,-75.2\ ,-2,75.2\ ,75.2\ ,4]$ | $[0.1\ ,0.1\ ,0.15]$ | Adam | $1\times10^{-2}$ | $1\times10^{-3}$ | 32 | 30 |
| Tab. 6 | KITTI [12] | SECOND [61] | $[\ \ \ 0.0\ ,-40.0\ ,-3,70.4\ ,40.0\ ,1]$ | $[0.16,0.16,6.0\ ]$ | Adam | $3\times10^{-3}$ | $1\times10^{-2}$ | 32 | 80 |
| Tab. 6a | Waymo [50] | DSVT [55] | $[-74.88,-74.88,-2,74.88,74.88,4]$ | $[0.32,0.32,6.0\ ]$ | Adam | $3\times10^{-3}$ | $5\times10^{-2}$ | 24 | 30 |
| Tab. 6b | Waymo [50] | Voxel R-CNN [8] | $[-75.2\ ,-75.2\ ,-2,75.2\ ,75.2\ ,4]$ | $[0.1\ ,0.1\ ,0.15]$ | Adam | $1\times10^{-2}$ | $1\times10^{-3}$ | 32 | 30 |

Table 8. Complete experimental setup for each table from the main manuscript. We specify the source domain, where $*$ specifies all except the target dataset for our multi-source domain generalization. We report LiDAR point cloud range, voxel size, optimizer parameters (learning rate (LR), weight decay (WD)), batch size (BS) and the number of epochs (E) used for training.

| $z$-alignment | Method | Car | Pedestrian | Cyclist | mAP |
|---|---|---|---|---|---|
| ✓ | Voxel R-CNN [8] | 66.93/28.80 | 23.39/18.65 | 19.23/15.76 | 36.52/21.07 |
| | Voxel R-CNN [8] w/ GBlobs | **80.95/53.98** | **38.33/33.22** | **29.18/25.68** | **49.48/37.62** |
| ✗ | Voxel R-CNN [8] | 54.61/20.83 | 10.51/ 7.68 | 5.88/ 5.12 | 23.66/11.21 |
| | Voxel R-CNN [8] w/ GBlobs | **80.84/55.05** | **37.93/33.24** | **28.62/24.60** | **49.13/37.63** |

Table 9. Influence of $z$-alignment on detector trained with different input features. We trained Voxel R-CNN [8] on nuScenes [2] using all three classes (Car, Pedestrian, Cyclist) simultaneously and evaluated performance using Average Precision (AP) on Bird's-eye View (BEV) / 3D views at 40 recall positions. Intersection over Union (IoU) thresholds of 0.7, 0.5, and 0.5 were used for Car, Pedestrian, and Cyclist, respectively. We evaluate the performance on KITTI [12], where we report the average AP across all difficulty levels (Easy, Moderate, Hard). Additionally, we provide the mean AP over the three classes. The best value in each category is highlighted in bold.
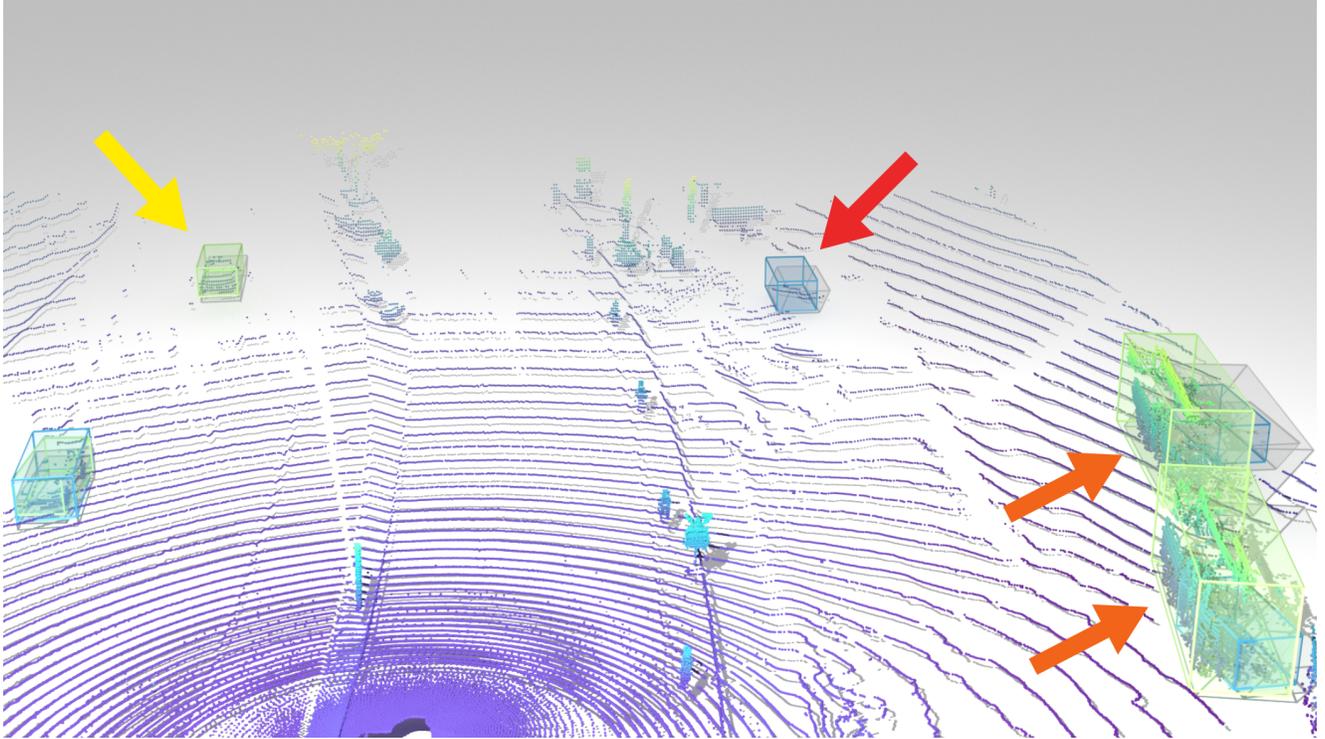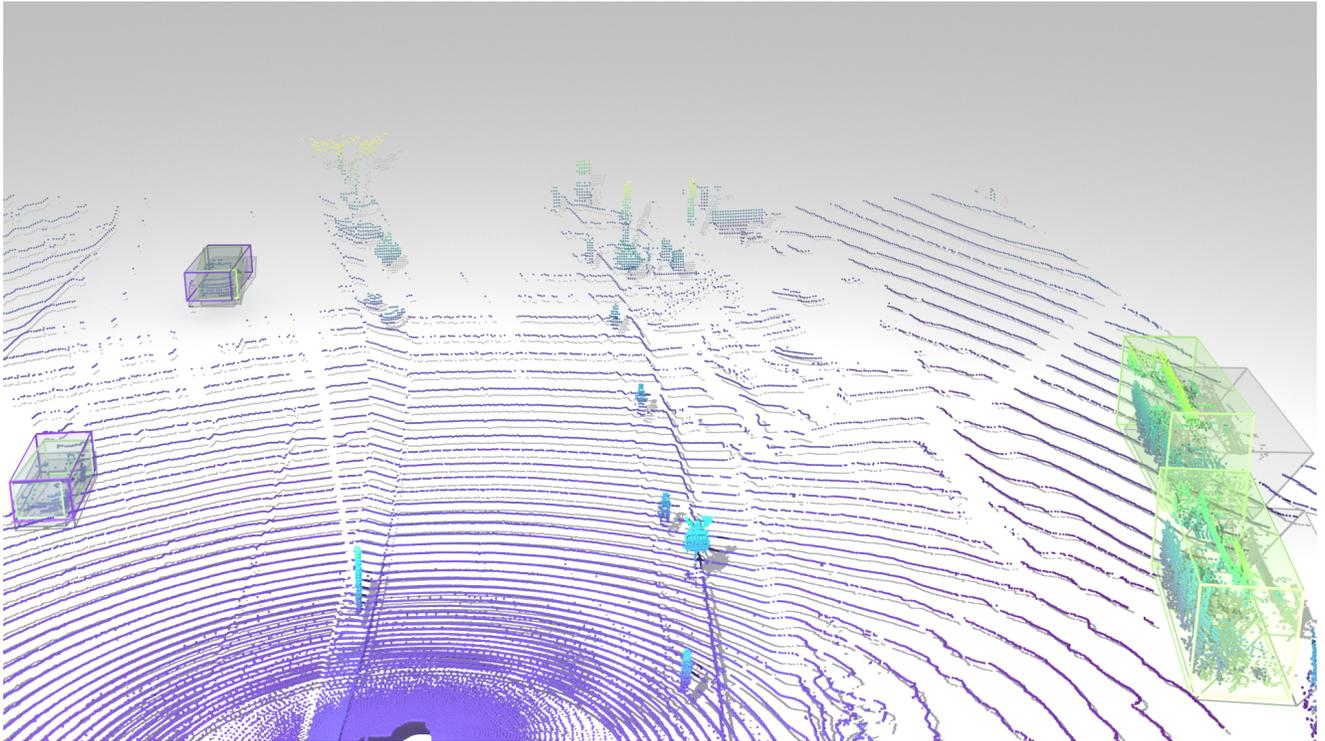
(a) nuScenes→KITTI Voxel R-CNN [8].



(b) nuScenes→KITTI Voxel R-CNN [8] w/ GBlobs.

Figure 6. Qualitative evaluation of Voxel R-CNN [8] on a nuScenes→KITTI benchmark thresholded at 0.5. Ground truth detections are shown in green. Detections from a model trained on standard global input features and our GBlobs are depicted in blue (a) and purple (b), respectively. False positive detections are marked with red arrows. The color of the point cloud represents the height.

(a) KITTI→Waymo SECOND [61].



(b) KITTI→Waymo SECOND [61] w/ GBlobs.

Figure 7. Qualitative evaluation of SECOND [61] on a KITTI→Waymo benchmark thresholded at 0.5. Ground truth detections are shown in green. Detections from a model trained on standard global input features and our GBlobs are depicted in blue (a) and purple (b), respectively. False positive and false negative detections are marked with red and yellow arrow, respetively. Detections which are dubious are markes with orange arrow. The color of the point cloud represents the height.