# LidarAugment: Searching for Scalable 3D LiDAR Data Augmentations

Zhaoqi Leng[1*], Guowang Li[1], Chenxi Liu[1], Ekin Dogus Cubuk[2], Pei Sun[1], Tong He[1],
Dragomir Anguelov[1] and Mingxing Tan[1]

*Abstract*— Data augmentations are important in training high-performance 3D object detectors for point clouds. Despite recent efforts on designing new data augmentations, perhaps surprisingly, most state-of-the-art 3D detectors only use a few simple data augmentations. In particular, different from 2D image data augmentations, 3D data augmentations need to account for different representations of input data and require being customized for different models, which introduces significant overhead. In this paper, we resort to a search-based approach, and propose *LidarAugment*, a practical and effective data augmentation strategy for 3D object detection. Unlike previous approaches where all augmentation policies are tuned in an exponentially large search space, we propose to factorize and align the search space of each data augmentation, which cuts down the 20+ hyperparameters to 2, and significantly reduces the search complexity. We show LidarAugment can be customized for different model architectures with different input representations by a simple 2D grid search, and consistently improve both convolution-based UPillars/StarNet/RSN and transformer-based SWFormer. Furthermore, LidarAugment mitigates overfitting and allows us to scale up 3D detectors to much larger capacity. In particular, by combining with latest 3D detectors, our LidarAugment achieves a new state-of-the-art 74.8 mAPH L2 on Waymo Open Dataset.

## I. INTRODUCTION

Data augmentations are widely used in training deep neural networks. In particular, for autonomous driving, many data augmentations are developed to improve data efficiency and model generalization. However, most recent 3D object detectors only use a few basic data augmentation operations such as rotation, flip and ground-truth sampling [1], [2], [3], [4], [5], [6], [7]. This is in a surprising contrast to 2D image recognition and detection, where much more sophisticated 2D data augmentations are commonly used in modern image-based models [8], [9], [10], [11], [12], [13]. In this paper, we aim to answer: *is it practical to adopt more advanced 3D data augmentations to improve modern 3D object detectors, especially for high-capacity models?*

The main challenge of adopting advanced 3D data augmentations is that 3D augmentations are often sensitive to input representations and model capacity. For example, range image based models and point cloud based models require different types of data augmentation due to different input representations. High capacity 3D detectors are typically prone to overfitting and require stronger overall data augmentation compared to lite models with fewer parameters. Therefore, tailoring each 3D augmentation for different models is necessary. However, the search space scales exponentially with respect to the number of hyperparameters, which leads
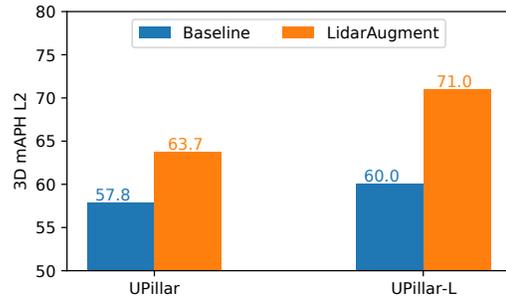
Fig. 1: **Model scaling with LidarAugment on Waymo Open Dataset.** Baseline augmentations are from the prior art of [14]. When scaling up UPillars to UPillars-L, our LidarAugment improves both models, and the gains are more significant for the larger model, thanks to its customizable regularization. More results in Table IV.

to significant search cost. Recent studies [15], [16] attempt to address these challenges by using efficient search algorithms. Those approaches typically construct a fixed search space, and run a complex search algorithms (such as population-based search [17]) to find a data augmentation strategy for a model. However, our studies reveal that the search spaces used in prior works are suboptimal. Despite having complex search algorithms, without a systematic way to define a good search space, we cannot unleash the potential of a model.

In this paper, we propose LidarAugment, a simplified search-based approach for 3D data augmentations. Unlike previous methods that rely on complex search algorithms to explore an exponentially large search space, our approach aims to define a simplified search space that contains a variety of data augmentations but has minimal (i.e. two) hyperparameters, such that users can easily customize a diverse set of 3D data augmentations for different models.

Specifically, we construct the LidarAugment search space by first factorizing a large search space based on operations and exploring each sub search space with a per-operation search. Then, we normalize and align the sub search space for each data augmentation to form the LidarAugment search space. The final LidarAugment search space contains only two shared hyperparameters: $m \in [0, \infty)$ controls the normalized magnitude and $p \in [0, 1]$ controls the probability of applying each data augmentation policies. Our LidarAugment search space significantly simplifies prior works [15] by cutting down the number of hyperparameters to two, a $15\times$ reduction in number of hyperparameters.

Despite only having two hyperparamters, our LidarAugment search space contains a variety of existing 3D data
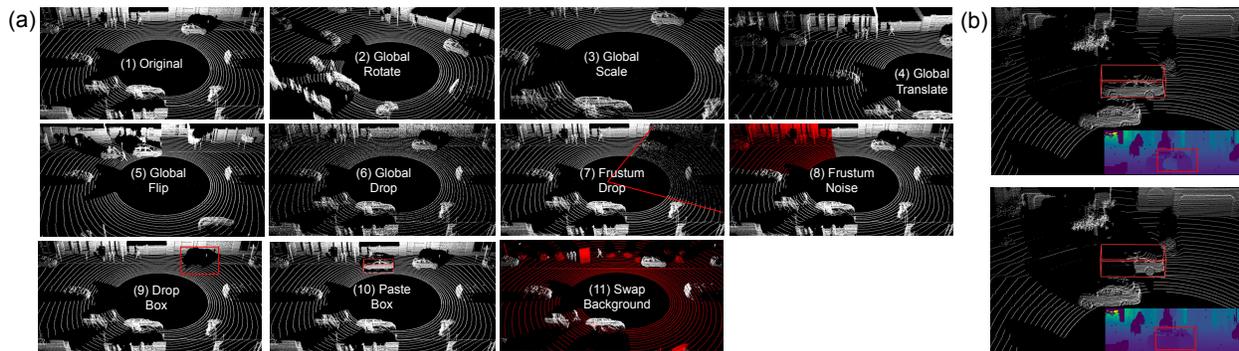
Fig. 2: **Visualizing LidarAugment.** (a) all data augmentation operations used in LidarAugment. For non-global operations, we highlight the augmented parts in red (boxes). (b) occlusion introduced by data augmentation, e.g., paste a car object, is handled by removing overlapping rays in range view based on distance. We show point clouds and the corresponding range images with (bottom)/without (top) removing overlapping rays in the range view.

augmentations, such as drop/paste 3D bounding boxes, rotate/scale/dropping points, and copy-paste objects and backgrounds. In addition, LidarAugment supports coherent augmentation across both point and range view representations, which generalizes to multi-view 3D detectors.

We perform extensive experiments on the Waymo Open Dataset [18] and demonstrate LidarAugment is effective and generalizes well to different model architectures (convolutions-based and transformer-based), different input views (3D point view and range image), and different temporal scales (single and multi frames). Notably, LidarAugment advances state-of-the-art (SOTA) transformer-based SWFormer by 1.4 mAPH on the test set. Furthermore, LidarAugment provides customizable regularization, which allows us to scale up 3D object detectors to much higher capacity without overfitting. As summarized in Figure 1, LidarAugment consistently improves UPillars models, and the performance gains are particularly large for high-capacity models. Our contributions can be summarized as:

1) **New insight:** we reveal that common 3D data augmentation search spaces are suboptimal and should be tailored for different models.
2) **LidarAugment:** we propose the LidarAugment search space, which supports jointly optimizing 10 augmentation policies with only two hyperparameters ($15\times$ reduction compares to prior works), offering diverse yet practical augmentations. In addition, we develop a new method to coherently augment both point and range-view input representations.
3) **State-of-the-art performance:** LidarAugment consistently improves both convolution-based UPillars/StarNet/RSN and attention-based SWFormer. With LidarAugment, we achieve new state-of-the-art results on Waymo Open Dataset. In addition, LidarAugment enables model scaling to achieve much better quality for high-capacity 3D detectors.

## II. RELATED WORKS

**Data augmentation.** Data augmentation is widely used in training deep neural networks In particular, for 3D object detection from point clouds, several global and local data augmentations, such as rotation, flip, pasting objects, and frustum noise, are used to improve model performance [19], [1], [20], [2], [4], [21], [15], [22], [23], [24]. However, as 3D data augmentations are sensitive to model architectures and capacity, it often requires extensive manual tuning to use these augmentations. Therefore, most existing 3D object detectors [2], [6], [25], [26], [14] only adopt a few simple augmentations, such as flip and shift pixels.

Several recent works attempt to use range images for multi-view 3D detection, but very few augmentations are developed for range images. [5] attempts to paste objects in the range image without handling occlusions. Our Paste Box augmentation support coherently augmenting both range-view and point-view input data while handling occluded objects in a simple way (more details in Figure 2), which enables more realistic augmented scenes and enriches the data augmentations for multi-view 3D detectors.

**Learning data augmentation policies.** Designing good data augmentation normally requires manual tuning and domain expertise. Several search-based approaches have been proposed for 2D images, such as AutoAugment [9], RandAugment [12], and Fast AutoAugment [27]. Our LidarAugment is inspired by RandAugment in the sense that we both try to construct a simplified search space. However, unlike 2D image augmentations, where a search space works well for many models, we reveal that existing search space for 3D detection tasks are suboptimal, which motivates us to propose the first systematical method to define search spaces for 3D detection tasks.

On the other hand, for 3D detection, PPBA [15] and PointAugment [16] propose efficient learning-based data augmentation frameworks for 3D point clouds. However, both works require users to run a complex algorithm on an exponentially large but not well-designed search space. In contrast, our work provides a systematical framework to

design a simple and more effective search spaces with only two hyperparameters.

## III. LidarAugment

In this section, we first introduce data augmentation policies used in LidarAugment. Next, we analyze the performance of each data augmentation policy on Waymo Open Dataset [18]. Finally, we propose a systematic approach to progressively design 3D augmentation search space.

### A. Data augmentations for point clouds and range images.

3D point cloud and 2D range image are two different representations of LiDAR data. Despite being the native representation of LiDAR data, data augmentations for range image is not well studied compared to point clouds. Here, we revisit data augmentations for point clouds, and introduce a new method for coherently applying data augmentation to both point clouds and range images.

**Augmenting point clouds.** We follow the implementation of data augmentation policies described in recent studies [1], [15], [28], which contain global operations (rotate, scale, translate, flip, and drop points) and local operations (drop boxes, paste boxes, swap background, drop points and add feature noise in a frustum), shown in Figure 2 (a).

**Augmenting range images.** Different from sparse 3D point representation, pixels in range image are compact. Data augmentations, such as pasting objects and swap background, disturb the compact structure of range representation. Here, we propose a novel approach to coherently augment both 3D point view and 2D range view by leveraging the bijective property between point clouds and range images, while account for occlusion.

First, we transform the range image pixels to point cloud based on $(x, y, z)$ coordinates. To preserve the bijective mapping between a pixel in a range image and a point in the corresponding point clouds, we concatenate the (row, column) index of each pixel in the range image as additional features before scattering pixels to 3D. After performing data augmentation in the point representation, we transform the augmented point clouds back to the range view by scattering each point to a pixel in a 2D image based on its (row, column) index.

**Leveraging the compactness of range images.** Coherently augmenting both range and point views leads to more realistic augmented scenes. Because each pixel in a range image corresponds to a unique ray from LiDAR, overlapping pixels in the range view represent that the same light ray penetrates trough multiple surfaces. When this happens, we compare the distance among overlapping pixels in the range view and keep the pixel that is closest to the ego vehicle. This effectively removes occluded points in both the range and point views, as shown in Figure 2 (b).

### B. Effects of each data augmentation.

In this section, we assess the effects of each data augmentation policy on Waymo Open Dataset [18]. To benchmark the policies, we develop a UPillars architecture, which is

| Policy | Hyperparameters | WOD (Veh./Ped.) | mAP L1 |
|---|---|---|---|
| No Aug | - | - | 60.2 |
| Drop Box | Probability<br>Number of boxes | $p/p$<br>$2m/2.8m$ | 66.0 (+5.8) |
| Paste Box | Probability<br>Number of boxes | $1.4p/p$<br>$3.2m/4.4m$ | 66.6 (+6.4) |
| Swap Background | Probability | $0.6p$ | 63.6 (+3.4) |
| Global Rot | Probability<br>Max rotation angle | $1.4p$<br>$0.22\pi m$ | 73.3 (+13.1) |
| Global Scale | Probability<br>Scaling factor | $p$<br>$0.036m$ | 66.0 (+5.8) |
| Global Drop | Probability<br>Drop ratio | $p$<br>$1 - 0.18m$ | 64.9 (+4.7) |
| Frustum Drop | Probability<br>Theta angle width<br>Phi angle width<br>R distance<br>Drop ratio | $p$<br>$0.1\pi m$<br>$0.1\pi m$<br>$75 - 7.5m$<br>$1 - 0.1m$ | 64.1 (+3.9) |
| Frustum Noise | Probability<br>Theta angle width<br>Phi angle width<br>R distance<br>Max noise level | $0.6p$<br>$0.14\pi m$<br>$0.14\pi m$<br>$75 - 10.5m$<br>$0.14m$ | 65.1 (+4.9) |
| Global Translate | Probability<br>Stdev. of noise (x, y) | $1.4p$<br>$0.66m$ | 67.5 (+7.3) |
| Global Flip | Probability | $p$ | 69.0 (+8.8) |

TABLE I: **Aligned search spaces and performance.** The search space of each hyperparameter for Waymo Open Dataset (WOD) for UPillars is listed. $(p, m)$ are two global hyperparameters to control all data augmentation policies. After aliging the search space, the optimal $(p, m)$ for each data augmentation are $(0.5, 5)$. The probability of each policy is clipped to [0, 1]. The min R distance is clipped to 0. The maximum rotation angle is clipped to $[0, \pi]$. The maximum flip probability is clipped to 0.5. The ratio of dropped points are clipped to [0, 0.8]. The theta angle and phi angle are clipped to $[0, \pi]$ and $[0, 2\pi]$, respectively.

based on the popular PointPillars [2], but incorporates recent optimizations in architecture design, i.e., unet backbone [29], and center net detection head [30].

**Datasets and training.** Waymo Open Dataset [18] contains 798 and 202 training and validation sequences. For the following studies, we train UPillars with batch size 64, Adam optimizer [31] and cosine decay learning rate with max learning rate 3e-3 and total step 80000.

**Effect of each data augmentations.** We factorize the LidarAugment search space into per-policy sub search space and show the UPillars performance when trained using only one policy on Waymo Open Dataset (WOD) in Table I. Interestingly, on WOD, the most effective data augmentation technique is global rotation, whereas on KITTI [32], pasting ground truth bounding boxes is commonly regarded as the most effective data augmentation [1]. A closer look at the statistics of the two datasets reveals that, on average, each KITTI LiDAR frame contains about five objects, whereas, each frame in WOD on average contains more than 50 objects. Thus, pasting ground truth objects has larger impact on KITTI, due to the significantly lower object density, than WOD. On the other hand, smaller global rotation angle $\pi/4$ is commonly used when training KITTI dataset, but we find much stronger rotation $\pi$ is preferred for WOD.
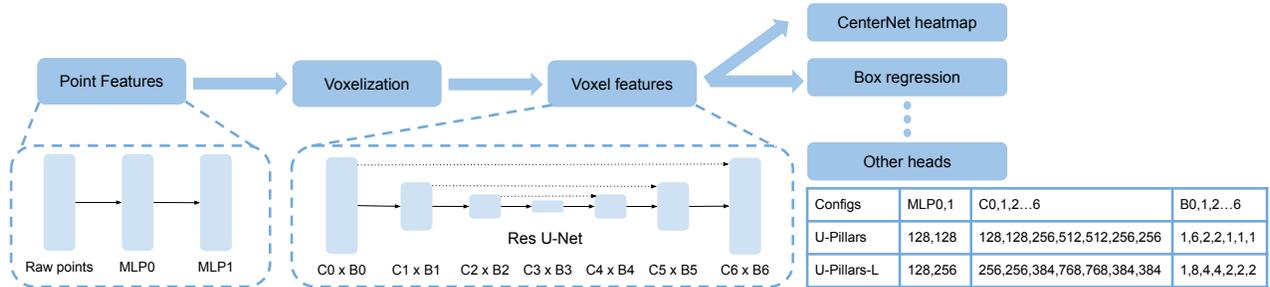
Fig. 3: **UPillars architecture.** Input points are processed by two full connected layers with channel size (MLP0, MLP1) before voxelized into pillars. The bird-eye-view pillars are processed by a Res U-Net, where the channel size and number of blocks at each resolution are (Ci, Bi). CenterNet detection head, box regression, and other attributes regression heads are applied to the output of U-Net. For both models, we use the same voxel size $0.32m$ and range $81.92m$.

## C. Defining LidarAugment search space.

As indicated from the previous section, different from RandAugment for 2D images [12], naively using the same search space across different datasets is suboptimal, which is a unique challenge for 3D detection tasks. To mitigate this new challenge, we propose to factorize the whole search space and align each data augmentation based on its optimal hyperparameters.

**Align the search space.** Using global hyperparameters to control all data augmentations requires normalizing the search domain of each hyperparameter. Without normalization, the same global magnitude could lead to an aggressive application of one data augmentation, and an insufficient application of another data augmentation. To align the search domain of each data augmentation policy, we train a UPillars model on a given dataset while only applying a single data augmentation policy at each time. Based on the optimal values of hyperparameters, we rescale the search domain of each hyperparameter in a data augmentation policy from $[0, arbitrary\_value]$ to $[0, optimal\_value]$ such that the optimal value for each hyperparameter corresponds to the same global magnitude or probability hyperparameter. Since each data augmentation policy contains multiple parameters, to save cost, we perform a small scale 2D grid search to scale the probability and magnitudes of all hyperparameters in each sub search space.

Here, we use Global Translate as an example. We define the initial search domain for the probability of applying Global Translate data augmentation to be $\{0.3, 0.5, 0.7, 0.9\}$, and the domain for the magnitude of translation noise

```
augmentations = [
    DropBox, PasteBox, SwapBackground, GlobalRot,
    GlobalScale, GlobalDrop, FrustumDrop ,
    FrustumNoise, GlobalTranslate, GlobalFlip]
def lidaraugment(m, p, input_frame):
    for aug in augmentations:
        aug.set_magnitude(m)
        aug.set_probability(p)
        input_frame = aug.transform(input_frame)
    return input_frame
```

Fig. 4: **Pseudo Python code for LidarAugment.**

to be $\{0.9, 1.5, 2.1, 2.7, 3.3, 3.9\}$. If the optimal values $(p_{\text{noise}}, m_{\text{noise}}) = (0.7, 3.3)$, we rescale the search domain to $(p_{\text{noise}}, m_{\text{noise}}) = (1.4p, 0.66m)$ such that when the global hyperparameters $(p, m) = (0.5, 5)$, the hyperparameters for Global Translate are optimal. Details about all the hyperparameters are listed in Table I. The LidarAugment pseudocode is shown in Figure 4.

## IV. EXPERIMENTS

In this section, we first introduce our experimental setups. Then we show LidarAugment significantly improves the performance for both convolution-based and attention-based models. Lastly, we show the model scaling results, followed by ablations studied on different models and datasets.

## A. Experimental setup

Our experiments are mostly based on Waymo Open Dataset [18] (WOD) where the main metric is mAPH L2, with additional ablation studies on nuScene [33]. We evaluate LidarAugment on a variety of 3D object detectors, as well as different model sizes. For fair comparison, we strictly follow the original training settings for each model, and only replace the baseline augmentation with our new LidarAugment. We train UPillars using Adam optimizer [31] and apply cosine learning rate with max learning rate 1e-3, total 16e4 steps and batch size 64.

## B. LidarAugment achieves new state-of-the-art results

Table II compares the validation set results on Waymo Open dataset. Our LidarAugment significantly improves both convolution-based and transformer-based models. In particular, by scaling up the basic UPillars, our LidarAugment achieves 71.0 mAPH of L2 on UPillars-L, which is 1.9 AP better than the previous best convolution-based 3D detector PVRCNN++ [34]. Notably, the latest transformer-based SWFormer [14] already uses 4 strong data augmentation policies, i.e., rotation (probability 0.74, yam angle uniformly sampled from $[-\pi, \pi]$), random flip (probability 0.5), randomly scaling the world (scaling factor uniformly sampled from $[0.95, 1.05]$, and randomly drop points (drop probability 0.05), where the rotation angle and flip probability are

| Method | Type | mAPH L2 | Vehicle AP/APH 3D L1 | L2 | Pedestrian AP/APH 3D L1 | L2 |
|---|---|---|---|---|---|---|
| P.Pillars [2] † | conv | 51.9 | 63.3/62.7 | 55.2/54.7 | 68.9/56.6 | 60.4/49.1 |
| CenterPoint [25] | conv | 67.1 | 76.6/76.1 | 68.9/68.4 | 79.0/73.4 | 71.0/65.8 |
| RSN_3f [6] | conv | 68.1 | 78.4/78.1 | 69.5/69.1 | 79.4/76.2 | 69.9/67.0 |
| PVRCNN++ [34] | conv | 69.1 | 79.3/78.8 | 70.6/70.2 | 81.8/76.3 | 73.2/68.0 |
| UPillars-L† | conv | 60.0 | 69.5/69.0 | 61.5/61.0 | 70.4/66.1 | 63.0/59.0 |
| **UPillars-L(+LA)** | **conv** | **71.0** | **79.5/79.0** | **71.9/71.5** | **81.5/77.3** | **74.5/70.5** |
| SST_1f [26] | attn | 63.4 | 74.2/73.8 | 65.5/65.1 | 78.7/69.6 | 70.0/61.7 |
| SST_3f [26] | attn | 69.5 | 77.0/76.6 | 68.5/68.1 | 82.4/78.0 | 75.1/70.9 |
| SWFormer [14] | attn | 70.9 | 79.4/78.9 | 71.1/70.6 | 82.9/79.0 | 74.8/71.1 |
| **SWFormer(+LA)** | **attn** | **72.8** | **80.9/80.4** | **72.8/72.4** | **84.4/80.7** | **76.8/73.2** |

TABLE II: **WOD validation-set results.** *LA* denotes our LidarAugment, *conv* denotes convolutional networks, and *attn* denotes attention-based transformer models. LidarAugment improves both types of models, and achieves the best results among each category. † model is trained using augmentations shown in prior art [14].

maxed out. Despite that, LidarAugment still outperforms SWFormer by 1.9 AP, establishing a new state-of-the-art result for single-modal models without ensemble or test time augmentation on Waymo Open Dataset.

Table III compares the test-set results among latest models. Compared to the latest SWFormer, our LidarAugment improves the test-set L2 mAPH by 1.4 AP, outperforming all prior arts by a large margin.

| Method | mAPH L2 | Vehicle AP/APH 3D L1 | L2 | Pedestrian AP/APH 3D L1 | L2 |
|---|---|---|---|---|---|
| P.Pillars [2] † | 55.1 | 68.6/68.1 | 60.5/60.1 | 68.0/55.5 | 61.4/50.1 |
| CenterPoint [25] | 69.1 | 80.2/79.7 | 72.2/71.8 | 78.3/72.1 | 72.2/66.4 |
| RSN_3f [6] | 69.7 | 80.7/80.3 | 71.9/71.6 | 78.9/75.6 | 70.7/67.8 |
| PVRCNN++ [34] | 71.2 | 81.6/81.2 | 73.9/73.5 | 80.4/75.0 | 74.1/69.0 |
| SST_TS_3f [26] | 72.9 | 81.0/80.6 | 73.1/72.7 | 83.1/79.4 | 76.7/73.1 |
| SWFormer [14] | 73.4 | 82.9/82.5 | 75.0/74.7 | 82.1/78.1 | 75.9/72.1 |
| **SWFormer(+LA)** | **74.8** | **84.0/83.6** | **76.3/76.0** | **83.1/79.3** | **77.2/73.5** |

TABLE III: **WOD test-set results.** LidarAugment (LA) significantly improves detection performance for SWFormer, and achieves new state-of-the-art mAPH L2.

### C. LidarAugment enables better model scaling

Scaling up model capacity is a common approach to achieve better performance, but large 3D object detectors often suffer from overfitting. Table IV shows scaling results, where UPillars-L is a larger model with more layers and channels than UPillars, detailed in Figure 3.

Here, we adopt the strong data augmentations used in the latest SWFormer as Baseline (see subsection IV-B). As shown in Table IV, with Baseline augmentations, UPillars-L does not benefit much from its significantly larger capacity. In fact, several metrics, such as Veh/Ped L1 AP, even become worse (e.g. 69.3/70.3 for UPillar-L vs 72.1/72.3 for UPillars). We observe the training loss of UPillars-L is much smaller compared to loss of UPillars, indicating severe overfitting.

On the other hand, LidarAugment achieves much better performance on larger models, especially on the most challenging metric, i.e., +7.3AP for 3D L2 mAPH as shown in Figure 1. Perhaps surprisingly, although the baseline UPillar (mAPH=57.8) is much worse than latest 3D detectors, the final performance of UPillar-L (+ LidarAugment) is actually

competitive with the latest SWFormers, i.e., their mAPH are 71.0 vs. 72.8. This opens up new research opportunities on exploring much larger and higher performance 3D detectors in the future.

| | Veh/Ped AP L1 UPillars | UPillars-L | Veh/Ped APH L2 UPillars | UPillars-L |
|---|---|---|---|---|
| BaseAugment | 72.1/72.3 | 69.3/70.3 | 63.5/52.1 | 61.0/59.0 |
| LidarAugment | **77.1/77.5** | **79.5/81.6** | **68.5/58.9** | **71.5/70.5** |

TABLE IV: **UPillars scaling results on WOD**.

### D. LidarAugment supports different representations

Different from 2D image models, 3D detectors are more diverse and could utilize different input representations due to the additional dimensionality and sparsity of point cloud data. Other than UPillars and SWFormer, which are both pillar-based architectures and taking 3D sparse points as inputs, we further demonstrate LidarAugment generalizes to other input representations. First, StarNet [35] is a point-based detector which directly processes raw points in 3D to detect objects. RSN, on the other hand, utilize multi-view property of point clouds and takes both range images and 3D sparse points as inputs. However, due to the lack of multi-view data augmentations in prior works, RSN only utilize two simple augmentations, i.e. random flip and rotation.

| Model | Augmentation | Vehicle | Pedestrian |
|---|---|---|---|
| StarNet [35] | baseline | 58.2 | 71.9 |
| | **+LidarAugment** | **61.6** | **74.2** |
| RSN-1frame [6] | baseline | 75.2 | 77.2 |
| | **+LidarAugment** | **75.8** | **79.0** |
| RSN-3frame [6] | baseline | 77.0 | 79.1 |
| | **+LidarAugment** | **77.7** | **80.6** |
| SWFormer [14] | baseline | 79.4 | 82.9 |
| | **+LidarAugment** | **80.9** | **84.4** |

TABLE V: **LidarAugment improves various models**. Start-net is a point-based detector. RSN is a range image and pillar-based detector. Results are WOD L1 AP.

LidarAugment is a general method which supports augmenting different views of point clouds, including range images, as explained in subsection III-A. Table V shows the performance of LidarAugment on point-based StarNet, range image based RSN, and transformer-based SWFormer. In general, our LidarAugment improves all kinds of 3D detectors, sometimes by a large margin.

### E. Ablation studies: comparing to other approaches.

In this section, we show LidarAugment outperforms other common data augmentation approaches on UPillars.

**Manually tuned data augmentation.** Due to the complexity of search space scales exponentially with respect to the number of parameters, commonly used data augmentation strategies often consists of few data augmentation operations. Here, we benchmark two sets of data augmentation strategies used in training high-performance 3D detectors. First, we adopt random flip (probability 0.5) and rotation (probability 0.5, yaw angle uniformly sampled from $[-\pi/4, \pi/4]$) data augmentations used in training RSN [6]. Then we benchmark more advanced and stronger data augmentation strategy used in training SWFormer [14], detailed in subsection IV-B. Our results show both data augmentation strategies significantly improved UPillars performances, about +10 AP for Vehicle and Pedestrian 3D L1 AP, when compared to the no augmentation baseline, shown in Table VI. However, tuning the data augmentation hyperparameters is challenging, e.g., if we only search 4 values for each hyperparameter, the number of searches of 5 hyperparameters exceeds 1000.

| UPillars (AP Level 1) | Hparams | Vehicle | Pedestrian |
|---|---|---|---|
| No Augmentation | - | 58.0 | 62.4 |
| Rotate & Flip [6] | 3 | 70.8 | 70.0 |
| Rotate, Flip, Scale, Drop points [14] | 5 | **72.1** | 72.3 |
| PPBA [15] | 29 | 71.6 | **72.6** |
| **LidarAugment** | **2** | **77.1** (+5.0) | **77.5** (+4.9) |

TABLE VI: **LidarAugment outperforms common data augmentation strategies.** UPillars L1 APs on Waymo Open Dataset *validation set* are reported. LidarAugment requires the least number of hyperparameters (Hparams) but achieves the best results compared to manually designed and automl-based data augmentations strategies.

**AutoML-based data augmentation.** To alleviate the challenge of exponentially large search space, population-based training is proposed to tune hyperparameters in data augmentations online [17], [11], [15]. We follow the implementation of progressive-population based data augmentation (PPBA) [15] and use the same sets of data augmentation policies and search space. We set population size 16, generation step 4000, perturbation and exploration rate to 0.2. Our results, in Table VI, show PPBA significant outperforms the no augmentation baseline. Despite PPBA introduces significantly more data augmentation policies, it is on par with manually tuned data augmentations, which only contains 4 policies.

We find the search space of PPBA is suboptimal after inspecting the search domain of each hyperparameter. For

example, the maximum rotation angle for global rotation in the PPBA search space is $\pi/4$, a common value used for KITTI dataset. However, $\pi/4$ is insufficient compared to the tailored max rotation angle $\pi$ used in our LidarAugment. Surprisingly, a single well-tuned global rotation augmentation achieves L1 mAP 73.3, shown in Table I, which outperforms PPBA with L1 mAP 72.1 over vehicle and pedestrian tasks. Although PPBA algorithm is more efficient than grid search and contains diverse augmentation policies, the suboptimal search domain of rotation angle restricts the performance of PPBA, which highlights the importance of tailoring 3D detection search space.

**LidarAugment** Alternatively, LidarAugment mitigates both the curse of dimensionality and suboptimal search domain issues by aligning and scaling the magnitude and probability of each data augmentation policy. This significantly reduces the search complexity (only 2 hyperparameters) while allowing exploration of a larger hyperparameter space. As indicate in Table VI, LidarAugment significantly outperforms both manually designed and AutoML-based data augmentation strategies by about 5 AP for both vehicle and pedestrian detection tasks and only requires a simple grid search of two hyperparameters.

### F. Generalize to nuScenes dataset

To further validate our method, we evaluate LidarAugment on a different dataset: nuScenes [33]. For simplicity, we adopt the same training settings as Waymo Open Dataset, but reduce the voxel size to 0.25 and the total training steps by half for faster training. We use the same baseline augmentation as SWFormer, and redefine LidarAugment search space for nuScenes following subsection III-C. Table VII shows LidarAugment is a general approach, which outperforms the baseline augmentation by a large margin on nuScenes.

| UPillars | mAP | NDS |
|---|---|---|
| Rotate, Flip, Scale, Drop points [14] | 40.6 | 48.2 |
| **LidarAugment** | **46.7** (+6.1) | **53.4** (+5.2) |

TABLE VII: **nuScenes validation-set results**.

### V. CONCLUSION

In this paper, we propose *LidarAugment*, a scalable and effective 3D augmentation approach for 3D object detection. Based on the insight that 3D data augmentations are sensitive to model architecture and capacity, we propose a simplified search space, which contains two hyperparameters to control a diverse set of augmentations. LidarAugment outperforms both manually tuned and existing search-based data augmentation strategies by a large margin. Extensive studies show that LidarAugment generalizes to convolution and attention-based architectures, as well as point-based and range-based input representations. More importantly, LidarAugment significantly simplifies the search process for 3D data augmentations and opens up exciting new research opportunities, such as model scaling in 3D detection. With LidarAugment, we demonstrate new state-of-the-art 3D detection results on the challenging Waymo Open Dataset.

## REFERENCES

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

[2] A. H. Lang, S. Vora, H. Caesar, L. Zhou, J. Yang, and O. Beijbom, "Pointpillars: Fast encoders for object detection from point clouds," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12 697–12 705.

[3] Y. Zhou, P. Sun, Y. Zhang, D. Anguelov, J. Gao, T. Ouyang, J. Guo, J. Ngiam, and V. Vasudevan, "End-to-end multi-view fusion for 3d object detection in lidar point clouds," in *Conference on Robot Learning*. PMLR, 2020, pp. 923–932.

[4] S. Shi, C. Guo, L. Jiang, Z. Wang, J. Shi, X. Wang, and H. Li, "Pv-rcnn: Point-voxel feature set abstraction for 3d object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 529–10 538.

[5] Z. Liang, M. Zhang, Z. Zhang, X. Zhao, and S. Pu, "Rangercnn: Towards fast and accurate 3d object detection with range image representation," *arXiv preprint arXiv:2009.00206*, 2020.

[6] P. Sun, W. Wang, Y. Chai, G. Elsayed, A. Bewley, X. Zhang, C. Sminchisescu, and D. Anguelov, "Rsn: Range sparse net for efficient, accurate lidar 3d object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 5725–5734.

[7] A. Bewley, P. Sun, T. Mensink, D. Anguelov, and C. Sminchisescu, "Range conditioned dilated convolutions for scale invariant 3d object detection," in *Conference on Robot Learning*. PMLR, 2021, pp. 627–641.

[8] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," *arXiv preprint arXiv:1710.09412*, 2017.

[9] E. D. Cubuk, B. Zoph, D. Mane, V. Vasudevan, and Q. V. Le, "Autoaugment: Learning augmentation policies from data," *arXiv preprint arXiv:1805.09501*, 2018.

[10] S. Yun, D. Han, S. J. Oh, S. Chun, J. Choe, and Y. Yoo, "Cutmix: Regularization strategy to train strong classifiers with localizable features," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 6023–6032.

[11] D. Ho, E. Liang, X. Chen, I. Stoica, and P. Abbeel, "Population based augmentation: Efficient learning of augmentation policy schedules," in *International Conference on Machine Learning*. PMLR, 2019, pp. 2731–2741.

[12] E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le, "Randaugment: Practical automated data augmentation with a reduced search space," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 702–703.

[13] Z. Zhong, L. Zheng, G. Kang, S. Li, and Y. Yang, "Random erasing data augmentation," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 07, 2020, pp. 13 001–13 008.

[14] P. Sun, M. Tan, W. Wang, C. Liu, F. Xia, Z. Leng, and D. Anguelov, "Swformer: Sparse window transformer for 3d object detection in point clouds," *European Conference on Computer Vision (ECCV)*, 2022.

[15] S. Cheng, Z. Leng, E. D. Cubuk, B. Zoph, C. Bai, J. Ngiam, Y. Song, B. Caine, V. Vasudevan, C. Li *et al.*, "Improving 3d object detection through progressive population based augmentation," in *European Conference on Computer Vision*. Springer, 2020, pp. 279–294.

[16] R. Li, X. Li, P.-A. Heng, and C.-W. Fu, "Pointaugment: an auto-augmentation framework for point cloud classification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 6378–6387.

[17] M. Jaderberg, V. Dalibard, S. Osindero, W. M. Czarnecki, J. Donahue, A. Razavi, O. Vinyals, T. Green, I. Dunning, K. Simonyan *et al.*, "Population based training of neural networks," *arXiv preprint arXiv:1711.09846*, 2017.

[18] P. Sun, H. Kretzschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine *et al.*, "Scalability in perception for autonomous driving: Waymo open dataset," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 2446–2454.

[19] X. Chen, H. Ma, J. Wan, B. Li, and T. Xia, "Multi-view 3d object detection network for autonomous driving," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2017, pp. 1907–1915.

[20] B. Yang, W. Luo, and R. Urtasun, "Pixor: Real-time 3d object detection from point clouds," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2018, pp. 7652–7660.

[21] Y. Chen, V. T. Hu, E. Gavves, T. Mensink, P. Mettes, P. Yang, and C. G. Snoek, "Pointmixup: Augmentation for point clouds," in *European Conference on Computer Vision*. Springer, 2020, pp. 330–345.

[22] J. S. Hu and S. L. Waslander, "Pattern-aware data augmentation for lidar 3d object detection," in *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*. IEEE, 2021, pp. 2703–2710.

[23] J. Choi, Y. Song, and N. Kwak, "Part-aware data augmentation for 3d object detection in point cloud," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 3391–3397.

[24] M. Reuse, M. Simon, and B. Sick, "About the ambiguity of data augmentation for 3d object detection in autonomous driving," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 979–987.

[25] T. Yin, X. Zhou, and P. Krahenbuhl, "Center-based 3d object detection and tracking," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 11 784–11 793.

[26] L. Fan, Z. Pang, T. Zhang, Y.-X. Wang, H. Zhao, F. Wang, N. Wang, and Z. Zhang, "Embracing single stride 3d object detector with sparse transformer," *arXiv preprint arXiv:2112.06375*, 2021.

[27] S. Lim, I. Kim, T. Kim, C. Kim, and S. Kim, "Fast autoaugment," *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[28] Z. Leng, S. Cheng, B. Caine, W. Wang, X. Zhang, S. Jonathon, M. Tan, and A. Dragomir, "Pseudoaugment: Learning to use unlabeled data for data augmentation in point clouds," in *European Conference on Computer Vision*. Springer, 2022, pp. 279–294.

[29] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.

[30] T. Yin, X. Zhou, and P. Krahenbuhl, "Center-based 3d object detection and tracking," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 11 784–11 793.

[31] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[32] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1231–1237, 2013.

[33] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nuscenes: A multimodal dataset for autonomous driving," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 11 621–11 631.

[34] S. Shi, L. Jiang, J. Deng, Z. Wang, C. Guo, J. Shi, X. Wang, and H. Li, "Pv-rcnn++: Point-voxel feature set abstraction with local vector representation for 3d object detection," *arXiv preprint arXiv:2102.00463*, 2021.

[35] J. Ngiam, B. Caine, W. Han, B. Yang, Y. Chai, P. Sun, Y. Zhou, X. Yi, O. Alsharif, P. Nguyen *et al.*, "Starnet: Targeted computation for object detection in point clouds," *arXiv preprint arXiv:1908.11069*, 2019.