

SA-BEV: Generating Semantic-Aware Bird’s-Eye-View Feature for Multi-view 3D Object Detection

Jinqing Zhang¹, Yanan Zhang¹, Qingjie Liu^{1,2,3*}, Yunhong Wang^{1,3}

¹State Key Laboratory of Virtual Reality Technology and Systems, Beihang University, Beijing, China

²Zhongguancun Laboratory, Beijing, China

³Hangzhou Innovation Institute, Beihang University, Hangzhou, China

{zhangjinqing, zhangyanan, qingjie.liu, yhwang}@buaa.edu.cn

Abstract

Recently, the pure camera-based Bird’s-Eye-View (BEV) perception provides a feasible solution for economical autonomous driving. However, the existing BEV-based multi-view 3D detectors generally transform all image features into BEV features, without considering the problem that the large proportion of background information may submerge the object information. In this paper, we propose Semantic-Aware BEV Pooling (SA-BEVPool), which can filter out background information according to the semantic segmentation of image features and transform image features into semantic-aware BEV features. Accordingly, we propose BEV-Paste, an effective data augmentation strategy that closely matches with semantic-aware BEV feature. In addition, we design a Multi-Scale Cross-Task (MSCT) head, which combines task-specific and cross-task information to predict depth distribution and semantic segmentation more accurately, further improving the quality of semantic-aware BEV feature. Finally, we integrate the above modules into a novel multi-view 3D object detection framework, namely SA-BEV. Experiments on nuScenes show that SA-BEV achieves state-of-the-art performance. Code has been available at <https://github.com/mengtan00/SA-BEV.git>.

1. Introduction

Camera and LiDAR are the two most commonly used sensors for 3D object detection, which is essential to autonomous driving systems. LiDAR-based methods [4, 12, 33, 39, 34, 35, 42] have attained excellent performance due to the accurate spatial structure information of point clouds, but the expensive LiDAR sensor reduces its universality. In contrast, camera-based methods [30, 29, 31, 18, 19] are relatively low-cost with plentiful semantic information, but are constrained by the lack of geometric depth cues.

*indicates the corresponding author.

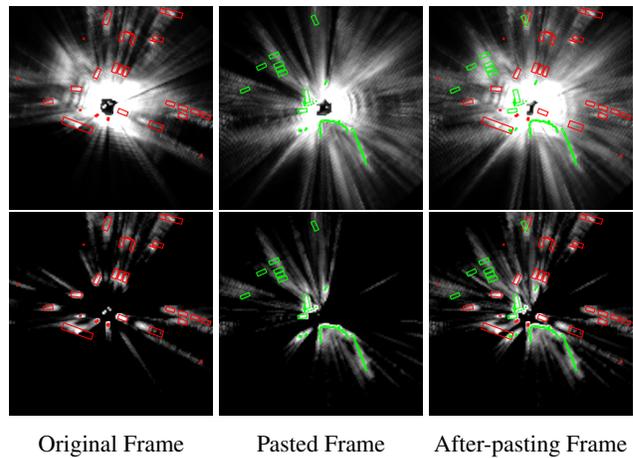


Figure 1: Comparison between normal BEV features (upper row) and semantic-aware BEV features (lower row). The brightness reveals the norm of the features and the red / green boxes are the ground truth of the original / pasted frame. The last column shows BEV-Paste, an data augmentation strategy that matches semantic-aware BEV features.

Considering the performance gap between camera and LiDAR, the Bird’s-Eye-View paradigm transforms multi-view image features into the BEV feature to make the following 3D perception easier [22, 10]. This practical and scalable camera-only paradigm is gaining popularity, and numerous advancements have allowed it to reach high perceptual precision [15, 14, 16, 8, 11]. The core step of the BEV paradigm is generating virtual points from image features, which will be projected into the “pillarized” BEV space. The features of the virtual points in the same pillar are then cumulated as the BEV feature. However, this operation does not fully utilize the semantic information of the image features and will inject massive background information that submerges object information.

In order to take full advantage of the valuable semantic

information of image features, we propose Semantic-Aware BEV Pooling (SA-BEVPool) to generate semantic-aware BEV features, which replace the normal BEV feature for 3D detection. Before projecting virtual points into BEV space, the semantic segmentation of image features is first predicted. If a virtual point is generated by the image element that belongs to the background, it will not be projected into BEV space. Similarly, virtual points with low depth scores will also be ignored. The comparison between normal BEV features and semantic-aware BEV features is shown in Fig. 1. SA-BEVPool can obviously filter out most of the background BEV features and alleviate the problem that the large proportion of background information submerges object information, therefore effectively improving the detection performance. Some multi-modal 3D object detectors [27, 36] also adopt segmentation on images when combining with LiDAR features, but they generally use powerful instance segmentation networks like CenterNet2 [43] to predict the segmentation of the large-scale image. Instead, SA-BEVPool can be easily applied in current BEV-based detectors like BEVDepth [15] and BEVStereo [14] by using their depth branch to simultaneously predict the semantic segmentation of small-scale image features.

GT-Paste [33] is a successful data augmentation strategy that has been frequently adopted by various LiDAR-based 3D detectors. However, due to the modality gap, it cannot directly adapt to camera-based 3D detectors. In our work, thanks to the reliable depth distribution and semantic segmentation predicated on image features, the semantic-aware BEV feature can approximately represent the information of all objects that are located appropriately in BEV space. As a result, adding the semantic-aware BEV features of another frame to the current semantic-aware BEV feature is the same as pasting all objects of another frame into the current frame. This strategy, we called BEV-Paste, enhances data diversity in a similar way to GT-Paste.

Although it is convenient to predict depth distribution and semantic segmentation with the same branch, doing so may result in a sub-optimal semantic-aware BEV feature. Research conclusion in the field of multi-task learning demonstrates that the integration of specific tasks and cross-task information is more conducive to the optimal solution of multiple prediction tasks. Inspired by this, we design a Multi-Scale Cross-Task (MSCT) head to combine the task-specific and cross-task information through multi-task distillation and dual-supervision on multiple scales prediction.

We integrate our proposed modules as a whole and name it SA-BEV. Extensive experiments on nuScenes dataset show that SA-BEV achieves a new state-of-the-art. In summary, the major contributions of this paper are:

- We propose SA-BEVPool, which uses semantic information to filter out unnecessary virtual points and generate the semantic-aware BEV feature, alleviating the

problem that the large proportion of background information submerges the object information.

- We propose BEV-Paste, an effective and convenient data augmentation strategy closely matching the semantic-aware BEV feature, which enhances data diversity and further promotes detection performance.
- We propose the MSCT head that combines the task-specific and cross-task information through multi-task learning on multiple scales, facilitating the optimization of the semantic-aware BEV feature.

2. Related Work

2.1. Vision-based 3D Object Detection

Although camera does not provide reliable depth of the surroundings like LiDAR, the plentiful semantic information carried by images still supports vision-based 3D object detectors to achieve considerable precision. Early vision-based 3D detectors predict attributes of 3D objects directly from 2D image features. For instance, CenterNet [44], a 2D detector, can be used to predict 3D objects without many modifications. Lately, FCOS3D [30] detects the 2D centers of the 3D objects, and features around the centers are used to predict the 3D attributes. PGD [29] establishes geometric relation graphs to improve the depth estimation results for better 3D object detection. DETR3D [31] follows DETR [2] and detects 3D objects with Transformer. PETR [18] introduces 3D position-aware representations, ameliorating the detection precision. PETRv2 [19] further brings in temporal information and improves efficiency.

Recently, an approach that transforms the image feature into the BEV feature is proposed by LSS [15], and BEVDet [10] employs the detection head of CenterPoint [35] to predict 3D objects from BEV feature. This paradigm can achieve comparable accuracy without cumbersome operations and is easy to extend, making it gain popularity. BEVDet4D [8] processes multiple key frames to introduce temporal information. BEVFormer [16] utilizes the deformable attention mechanism to generate the BEV feature. BEVDepth [15] applies explicit depth supervision on the predicted latent depth distribution, improving the detection accuracy. BEVStereo [14] further improves the quality of the depth by applying the multi-view stereo on nearby key frames. PolarFormer [11] generates the BEV feature using the polar coordinate for a more accurate location. However, these methods project all image features into BEV features, without considering the problem that the large proportion of background information may submerge the object information. In this paper, we propose SA-BEVPool, which can filter out background information according to the semantic segmentation of image features and generate semantic-aware BEV features.

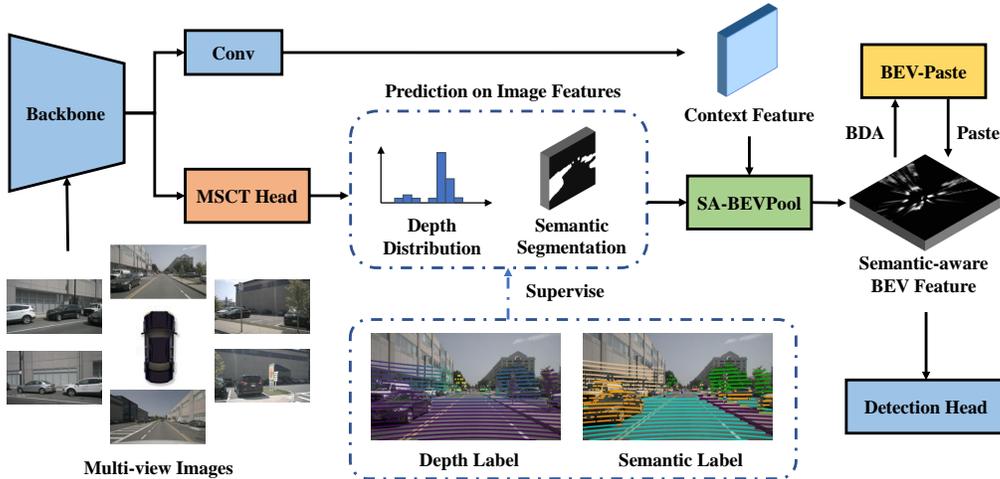


Figure 2: Overall framework of SA-BEV. The MSCT head uses multi-scale image features to predict the depth distribution and semantic segmentation, which are utilized by SA-BEVPool to generate the semantic-aware BEV feature. BEV-Paste is then applied to increase the diversity of BEV features during the training phase.

2.2. Data Augmentation in 3D Object Detection

The diversity of the dataset is crucial to the generalization performance of models. Besides regular data augmentation like random scaling, flipping and rotation, GT-Paste [33] is another effective strategy frequently used by LiDAR-based detectors. It crops the points according to the 3D boxes of ground truth and pastes them to other frames to create new training data. Lately, an improvement in generating a visibility map to correct the wrong occlusion relationship introduced by GT-Paste is proposed in [7]. The augmentation on the individual object is also proposed in [3, 41] which takes object points into parts and applies operations like dropout, swap, and mix on it.

Since GT-Paste shows excellent effectiveness on increase data diversity, there have been some attempts to adopt it in camera-only 3D detectors. Box-Mixup and Box-Cut-Paste proposed by [24] directly cut the objects from images according to their 2D bounding boxes and paste them into other frames. To paste precisely, objects are cropped by their instance masks in [37]. Pointaugmenting [28] utilizes a more complicated way to tackle the occlusion relationship of original objects and pasted objects. However, these attempts to expand the GT-Paste into image space cannot easily overcome the issues caused by the gap between LiDAR and camera. In this paper, we propose BEV-Paste, a convenient way to effectively extend GT-Paste into BEV-based methods with the help of SA-BEVPool.

2.3. Multi-task Learning

Multi-task learning generally leads to better prediction through interactive learning between multiple tasks. Ac-

cording to [25], both task-specific information and cross-task information are important for getting optimal results on multiple tasks. PAD-Net [32] proposes multi-modal distillation module to automatically supplement cross-task information. PAP-Net [40] extracts cross-task affinity patterns and recursively propagates the pattern by affinity matrices. MTI-Net [26] models the task interactions at different scales and aggregates multi-scale information to make precise predictions.

Some methods also introduce multi-task learning into 3D object detection. MMF [17] deeply fuses the features of images and LiDAR through simultaneous supervision made on multiple tasks. Latent support surfaces are estimated in [23] to help improve the precision of 3D detection. A multi-task LiDAR network proposed by [5] makes predictions on 3D detection and road understanding that can complement each other. Some BEV-based 3D detectors [16, 19] also apply BEV segmentation to obtain better BEV representation. In this paper, we propose the MSCT head that combines the task-specific and cross-task information of multiple scales for depth estimation and semantic segmentation, facilitating the optimization of the semantic-aware BEV feature.

3. Method

In this work, we propose SA-BEV, a novel multi-view 3D object detection framework that generates semantic-aware BEV features for better detection performance. It contains the Semantic-Aware BEV Pooling (SA-BEVPool), the BEV-Paste data augmentation strategy and the Multi-Scale Cross-Task (MSCT) head. The overall framework of SA-BEV is shown in Fig. 2.

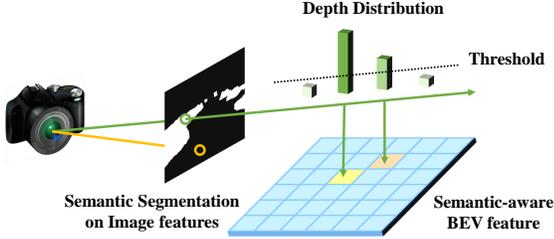


Figure 3: Illustration of Semantic-Aware BEV Pooling. The green line represents the projection of the foreground features, while the yellow line represents the ignored background features. The foreground virtual points with depth scores lower than the threshold are also ignored.

3.1. Semantic-Aware BEV Pooling

The way of transforming the image features into BEV features for better perception was first proposed by LSS [22]. It predicts the depth distribution α and context feature c of each image feature element. Then each element generates virtual points at different depths. The feature of the point at depth d is represented as $p_d = \alpha_d c$. After that, all virtual points will be projected to the BEV space which is divided into pillars. The features of virtual points in the same pillar will be cumulated as the BEV feature. This process is known as BEV pooling.

Subsequent BEV-based 3D detectors [15, 20, 9] significantly improve the efficiency and accuracy of BEV pooling. But what is unchanged is that these methods insist on projecting all virtual points into BEV space. However, we argue that this is unnecessary for 3D detection tasks. On the contrary, if all virtual points belonging to the background are projected, the foreground virtual points that account for less than 2% of the total virtual points will be submerged. It will confuse the following detection head and reduce the detection accuracy.

To highlight the valuable foreground information in the BEV features, we propose a novel Semantic-Aware BEV Pooling (SA-BEVPool) that is shown in Fig. 3. It applies semantic segmentation on the image features to get the foreground score β of each element. The element with low β is more possible to carry useless information for detection, and the virtual points generated from it will be ignored during the BEV pooling. Similarly, the virtual points with low α_d provide trivial information and will also be ignored. Denoting the filtering function as:

$$\mathcal{F}(x, y) = \begin{cases} 0, & x < y, \\ 1, & x \geq y, \end{cases} \quad (1)$$

the point features after filtering are changed as:

$$\hat{p}_d = \mathcal{F}(\alpha_d, T_D) \mathcal{F}(\beta, T_S) p_d, \quad (2)$$

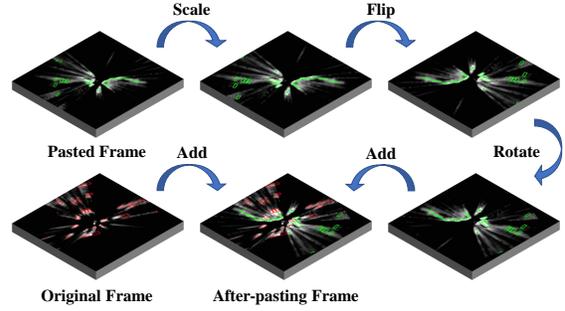


Figure 4: Illustration of BEV-Paste.

where T_D and T_S are the threshold for α_d and β . Only non-empty \hat{p}_d will construct BEV feature. Since the operation of filtering relatively low-value virtual points attaches semantic information to the generated BEV features, they can be called semantic-aware BEV features.

The difference between the normal BEV feature and the semantic-aware BEV feature is clearly shown in Fig. 1. The normal BEV features in the first row generally have a ring of light in the center, which represents the ground. It accounts for most of the signal strength in the feature without contributing useful information for detection. In contrast, most of the background information is removed in the semantic-aware BEV feature and object information is emphasized. Furthermore, the location of object information in the semantic-aware BEV feature matches the ground truth well, making the detection head easier to predict accurately.

3.2. BEV-Paste

GT-Paste [33] is a data augmentation strategy commonly used by LiDAR-based 3D detectors. It has been proved that the diversity of the dataset can be effectively increased by sampling the points in 3D boxes and pasting them into other frames. However, several problems prevent the application of GT-Paste in camera-based methods. First, sampling an object by the bounding box on the image can not get its pure data as the point cloud does. Another problem is that pasting objects to another image may wrongly occlude original objects and result in data loss. In addition, the illumination change of different frames also gives the pasted objects unnatural appearances. Some multi-modal 3D detectors [24, 28, 37, 38] make an effort to solve these issues but generally lack convenience and accuracy.

Here, we propose BEV-Paste that successfully applies GT-Paste in camera-only 3D detectors without complicated steps. With SA-BEVPool, the semantic-aware BEV features transformed from image features approximately represent the information of all objects in the frame as shown in Fig. 1. It makes adding arbitrary semantic-aware BEV features of two frames during the training phase equivalent to aggregating the objects contained in two frames into

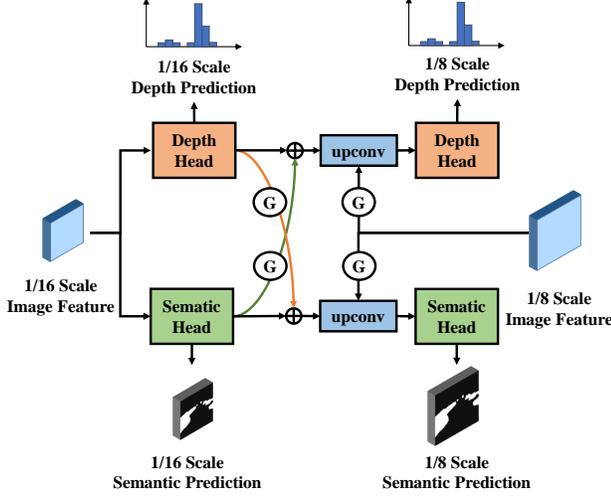


Figure 5: The structure of Multi-Scale Cross-Task head. Both 1/16 and 1/8 scale image features are taken as input.

one frame. While effectively increasing the diversity of the entire training dataset, BEV-Paste does not increase the computational cost in the inference stage.

In practice, we randomly select the original semantic-aware BEV feature \mathbf{B}_O and the pasted semantic-aware BEV feature \mathbf{B}_P from the same batch. This is to guarantee \mathbf{B}_O and \mathbf{B}_P follow the same distribution. Instead of directly pasting \mathbf{B}_P to \mathbf{B}_O , extra BEV data augmentations (BDA) shown in Fig. 4 are first applied to \mathbf{B}_P and $\hat{\mathbf{B}}_P$ is obtained. It prevents the data duplication of \mathbf{B}_P . The same augmentation is also applied to the ground truth of the pasted frame G_P to get \hat{G}_P . The detection loss after BEV-Paste can be represented as:

$$\mathcal{L}_{det} = \mathcal{L}_{det}(Det(\mathbf{B}_O + \hat{\mathbf{B}}_P), G_O \cup \hat{G}_P), \quad (3)$$

where Det includes BEV encoder and detection head, G_O is the ground truth of original frame.

3.3. Multi-Scale Cross-Task Head

It is a convenient way to obtain semantic-aware BEV features by making the depth branch predict semantic segmentation at the same time, but it generally leads to sub-optimal results. Let us regard the generation of the semantic-aware BEV feature as a multi-task learning application. According to the research conclusion, both task-specific information and cross-task information are important for getting the global optimal solution of multiple tasks. If the depth distribution and semantic segmentation are predicted by the same network branch, the network only extracts cross-task information from the image features and can not perform optimally on each task.

Inspired by the principle of multi-task learning, we design a Multi-Scale Cross-Task (MSCT) head as shown in

Fig. 5. In the first stage, the head takes 1/16 scale image feature \mathbf{F}_I^{16} as input and makes a relatively coarse prediction of depth distribution and semantic segmentation. After that, \mathbf{F}_I^{16} is transformed into depth feature \mathbf{F}_D^{16} and semantic feature \mathbf{F}_S^{16} , which carry the task-specific information of their own task. To complement cross-task information, Multi-Task Distillation (MTD) module proposed in [32] is applied between \mathbf{F}_D^{16} and \mathbf{F}_S^{16} . It is composed of several self-attention blocks, which generate gate map \mathcal{G} by

$$\mathcal{G}(\mathbf{F}) = \sigma(W_G \mathbf{F}), \quad (4)$$

where W_G is the gate convolution and σ denotes sigmoid function. The features supplemented by the cross-task information can be formulated as:

$$\hat{\mathbf{F}}_D^{16} = \mathbf{F}_D^{16} + \mathcal{G}(\mathbf{F}_D^{16}) \odot (W_t \mathbf{F}_S^{16}) \quad (5)$$

$$\hat{\mathbf{F}}_S^{16} = \mathbf{F}_S^{16} + \mathcal{G}(\mathbf{F}_S^{16}) \odot (W_t \mathbf{F}_D^{16}) \quad (6)$$

where W_t is the task convolution and \odot denotes element-wise multiplication. It is clear that MTD uses these self-attention blocks to automatically extract cross-task information from one task feature and add it to other task features.

After the task features interaction, $\hat{\mathbf{F}}_D^{16}$ and $\hat{\mathbf{F}}_S^{16}$ obtain both task-specific information and cross-task information. Before inputting them into the second stage prediction head, they are up-sampled to the 1/8 scale and combined with the 1/8 scale image feature \mathbf{F}_I^8 using the same self-attention blocks. The features can be formulated as:

$$\hat{\mathbf{F}}_D^8 = Up(\hat{\mathbf{F}}_D^{16}) + \mathcal{G}(\mathbf{F}_I^8) \odot (W_t \mathbf{F}_I^8) \quad (7)$$

$$\hat{\mathbf{F}}_S^8 = Up(\hat{\mathbf{F}}_S^{16}) + \mathcal{G}(\mathbf{F}_I^8) \odot (W_t \mathbf{F}_I^8) \quad (8)$$

The second stage head then predicts the relatively fine depth distribution and semantic segmentation which will be used to generate the semantic-aware BEV feature.

During training, both predictions on the 1/16 and 1/8 scale are supervised. It ensures the first stage head can extract task-specific information and the second stage head can combine the task-specific information with cross-task information. The supervision signals are obtained by projecting point clouds on images following BEVDepth [15]. The depth values of the projected points are the depth labels and the points in the 3D boxes are regarded as the foreground. The total loss can be formulated as:

$$\mathcal{L} = \mathcal{L}_{det} + \frac{\lambda_1}{2}(\mathcal{L}_S^{16} + \mathcal{L}_S^8) + \frac{\lambda_2}{2}(\mathcal{L}_D^{16} + \mathcal{L}_D^8) \quad (9)$$

4. Experiments

In this section, we first introduce our experimental settings. Then, comparisons with previous state-of-the-art multi-view 3D detectors are shown. Finally, comprehensive experiments with detailed ablation studies are conducted on SA-BEV to show the effectiveness of each component, i.e. SA-BEVPool, BEV-Paste and MSCT head.

Table 1: Comparison with previous state-of-the-art multi-view 3D detectors on the nuScenes *test* set.

| Method | Backbone | Resolution | mAP \uparrow | NDS \uparrow | mATE \downarrow | mASE \downarrow | mAOE \downarrow | mAVE \downarrow | mAAE \downarrow |
|------------------|------------|-------------------|----------------|----------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| FCOS3D [30] | ResNet-101 | 900 \times 1600 | 0.358 | 0.428 | 0.690 | 0.249 | 0.452 | 1.434 | 0.124 |
| PGD [29] | ResNet-101 | 900 \times 1600 | 0.386 | 0.448 | 0.626 | 0.245 | 0.451 | 1.509 | 0.127 |
| DETR3D [31] | V2-99 | 900 \times 1600 | 0.412 | 0.479 | 0.641 | 0.255 | 0.394 | 0.845 | 0.133 |
| BEVDet [10] | Swin-B | 900 \times 1600 | 0.424 | 0.488 | 0.524 | 0.242 | 0.373 | 0.950 | 0.148 |
| PETR [18] | V2-99 | 900 \times 1600 | 0.441 | 0.504 | 0.593 | 0.249 | 0.383 | 0.808 | 0.132 |
| BEVFormer [16] | V2-99 | 900 \times 1600 | 0.481 | 0.569 | 0.582 | 0.256 | 0.375 | 0.378 | 0.126 |
| BEVDet4D [8] | Swin-B | 640 \times 1600 | 0.451 | 0.569 | 0.511 | 0.241 | 0.386 | 0.301 | 0.121 |
| PolarFormer [11] | V2-99 | 900 \times 1600 | 0.493 | 0.572 | 0.556 | 0.256 | 0.364 | 0.440 | 0.127 |
| PETrv2 [19] | V2-99 | 640 \times 1600 | 0.490 | 0.582 | 0.561 | 0.243 | 0.361 | 0.343 | 0.120 |
| BEVDepth [15] | V2-99 | 640 \times 1600 | 0.503 | 0.600 | 0.445 | 0.245 | 0.378 | 0.320 | 0.126 |
| BEVStereo [14] | V2-99 | 640 \times 1600 | 0.525 | 0.610 | 0.431 | 0.246 | 0.358 | 0.357 | 0.138 |
| SA-BEV | V2-99 | 640 \times 1600 | 0.533 | 0.624 | 0.430 | 0.241 | 0.338 | 0.282 | 0.139 |

Table 2: Comparison with previous state-of-the-art multi-view 3D detectors on the nuScenes *val* set.

| Method | Backbone | Resolution | mAP \uparrow | NDS \uparrow |
|------------------|------------|-------------------|----------------|----------------|
| FCOS3D [30] | ResNet-101 | 900 \times 1600 | 0.343 | 0.415 |
| DETR3D [31] | ResNet-101 | 900 \times 1600 | 0.303 | 0.374 |
| PGD [29] | ResNet-101 | 900 \times 1600 | 0.369 | 0.428 |
| PETR [18] | ResNet-101 | 512 \times 1408 | 0.357 | 0.421 |
| BEVDet [10] | Swin-B | 900 \times 1600 | 0.393 | 0.472 |
| BEVFormer [16] | ResNet-101 | 900 \times 1600 | 0.416 | 0.517 |
| PETrv2 [19] | ResNet-101 | 900 \times 1600 | 0.421 | 0.524 |
| BEVDet4D [8] | Swin-B | 640 \times 1600 | 0.421 | 0.545 |
| PolarFormer [11] | ResNet-101 | 900 \times 1600 | 0.432 | 0.528 |
| BEVDepth [15] | ConvNeXt-B | 512 \times 1408 | 0.462 | 0.558 |
| BEVStereo [14] | ConvNeXt-B | 512 \times 1408 | 0.478 | 0.575 |
| SA-BEV | ConvNeXt-B | 512 \times 1408 | 0.479 | 0.579 |

4.1. Experimental Settings

4.1.1 Dataset and Metrics

nuScenes [1] dataset is a large-scale autonomous driving benchmark. It contains 750 scenarios for training, 150 scenarios for validation and 150 scenarios for testing. Each scenario lasts for around 20 seconds and the key samples are annotated at 2Hz. The data collected from six cameras, one LiDAR and five radars are provided to every sample. For 3D object detection, nuScenes Detection Score (NDS) is proposed to capture all aspects of the nuScenes detection tasks. Except mean average precision (mAP), NDS is also related to five types of true positive metrics (TP metrics), including mean Average Translation Error (mATE), mean Average Scale Error (mASE), mean Average Orientation Error (mAOE), mean Average Velocity Error (mAVE), mean Average Attribute Error (mAAE).

4.1.2 Implementation Detail

We accomplish our proposed improvements on the network structure of BEVDepth [15]. Our experiments are implemented based on MMDetection3D with 8 NVIDIA GeForce RTX 3090 GPUs. Models are trained with AdamW [21] optimizer and gradient clip is utilized. The universal data augmentation we adopt on the image and BEV feature follows the configuration in [10]. For the ablation study, we use ResNet-50 [6] as the image backbone and the image size is downsampled to 256 \times 704. The models are trained for 24 epochs without CBGS strategy [45] for the ablation study. When compared to other methods, the models are trained for 20 epochs with CBGS strategy.

4.2. Main Results

4.2.1 Comparison with State-of-the-Arts

We compare SA-BEV with state-of-the-art multi-view 3D detectors on nuScenes *test* set and show the results in Table 1. We take 640 \times 1600 resolution image as input and VoVNet-99 [13] as the image backbone. SA-BEV achieves the best mAP and NDS, 3.0% and 2.4% higher than its baseline (i.e. BEVDepth [15]). It also exceeds BEVStereo [14] by 0.8% mAP and 1.4% NDS, which adopts the complicated multi-view stereo structure for more accurate depth estimation. The comparison on nuScenes *val* set is shown in 2. It can be found SA-BEV also achieves the best detection precision. The decent results highlight the advantage of the proposed SA-BEV.

4.2.2 Visualization

We visualize the detection results on images and BEV features in Fig. 6. Compared to BEVDepth, SA-BEV can make more precise predictions with the help of semantic-aware BEV features. For instance, the orange dashed rect-

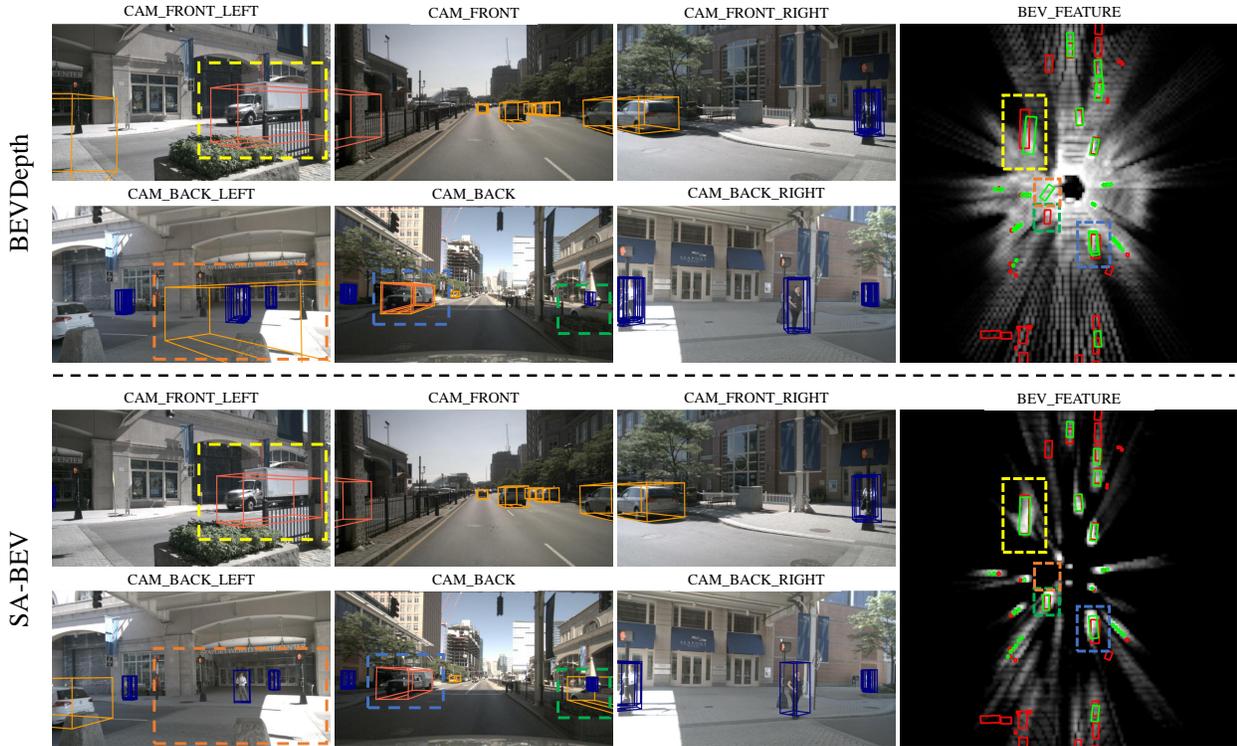


Figure 6: Visualization of detection results on images and BEV features. The red boxes and green boxes on BEV features represent the ground truth and the predicted boxes, respectively. The dashed rectangles illustrate that the prediction of SA-BEV is more precise than BEVDepth.

Table 3: Ablation study of component in SA-BEV on the nuScenes *val* set. *pool*, *paste* and *head* denotes SA-BEVPool, BEV-Paste and MSCT head, respectively.

| Baseline | <i>pool</i> | <i>paste</i> | <i>head</i> | mAP \uparrow | NDS |
|----------------|-------------|--------------|-------------|----------------|--------------|
| BEVDepth [15] | | | | 0.330 | 0.436 |
| | ✓ | | | 0.340 | 0.449 |
| | ✓ | ✓ | | 0.354 | 0.464 |
| | ✓ | ✓ | ✓ | 0.365 | 0.483 |
| BEVDet [10] | | | | 0.278 | 0.322 |
| | ✓ | ✓ | | 0.304 | 0.348 |
| BEVStereo [14] | | | | 0.349 | 0.454 |
| | ✓ | ✓ | | 0.364 | 0.467 |

angles show that the filtration of the background prevents SA-BEV from making false detection. The yellow dashed rectangles indicate that the semantic-aware BEV feature correctly emphasizes the location of the truck, which results in precise detection. In addition, the green / blue dashed rectangles display that SA-BEV can successfully recall the missed object and remove the redundant detection box.

4.3. Ablation Study

4.3.1 Component Analysis

We individually evaluate the contributions of SA-BEVPool, BEV-Paste and MSCT head with BEVDepth [15] as the baseline. The results are shown in Table 3. After applying SA-BEVPool, the performance is boosted by 1.0% and 1.3% on mAP and NDS. It is further improved by 1.4% / 1.5% and 1.1% / 1.9% through incorporating BEV-Paste and MSCT head respectively. Finally, we obtain the full model of SA-BEV, which gains 3.5% and 4.7% in total, validating its effectiveness. SA-BEVPool and BEV-Paste are also applied to BEVDet [10] and BEVStereo [15], increasing 2.6% / 2.6% and 1.5% / 1.3% respectively on mAP and NDS. It demonstrates that these components can be easily embedded into the existing BEV-based detectors and bring about noticeable precision improvement.

4.3.2 Semantic-Aware BEV Pooling

The semantic threshold T_S and semantic threshold T_D in SA-BEVPool control the scale of the valid virtual points. We vary thresholds and show the results in Table 4. The results indicate that even a low T_D can sharply reduce the

Table 4: Ablation study of the semantic threshold used in SA-BEVPool. “Percentage” denotes the average proportion of valid virtual points.

| T_D | T_S | mAP↑ | NDS↑ | Percentage |
|--------|-------|--------------|--------------|------------|
| - | - | 0.330 | 0.436 | 100% |
| 0.0085 | - | 0.338 | 0.438 | 7.92% |
| 0.0085 | 0.10 | 0.339 | 0.444 | 3.26% |
| 0.0085 | 0.25 | 0.340 | 0.449 | 1.80% |
| 0.0085 | 0.50 | 0.329 | 0.432 | 0.89% |

Table 5: Ablation study of BEV-Paste strategy. N_P denotes the average number of frames that are pasted to the original frame.

| Method | N_P | mAP↑ | NDS↑ |
|---------------|-------|--------------|--------------|
| w/o extra BDA | 0 | 0.340 | 0.449 |
| | 0.5 | 0.348 | 0.453 |
| | 1 | 0.349 | 0.453 |
| | 2 | 0.349 | 0.452 |
| w/ extra BDA | 1 | 0.354 | 0.464 |

Table 6: Ablation study of MSCT head. “MTD”, “DS” and “MS” denote the multi-task distillation module, the dual supervision and the utilization of multi-scale image features.

| MTD | DS | MS | mAP↑ | NDS↑ |
|-----|----|----|--------------|--------------|
| | | | 0.354 | 0.464 |
| ✓ | | | 0.358 | 0.468 |
| ✓ | ✓ | | 0.361 | 0.473 |
| | | ✓ | 0.361 | 0.478 |
| ✓ | ✓ | ✓ | 0.365 | 0.483 |

scale of valid virtual points and an appropriate T_S can effectively improve the detection precision. However, a too-high semantic threshold may lead to the loss of foreground information and damage the precision. We set the semantic threshold to 0.25 as a good trade-off. It only needs 1.8% valid virtual points, resulting in 1% mAP and 1.3% NDS performance improvement.

4.3.3 BEV-Paste

We conduct experiments with different settings when applying BEV-Paste, including the number of frames pasted to each original frame and whether to utilize extra BDA. The results are shown in Table 5. Setting N_P to 0.5 means half of the original frames are augmented by one pasted frame while the others are not augmented. As shown in Table 5, the BEV-Paste is not sensitive to the number of pasted

frames. Considering that pasting too many frames will cost more time on training detection head because the ground truth objects are increased, setting N_P as 1 is enough. Besides, extra BDA can effectively alleviate data duplication and further improve detection performance. The cooperation of these two points contributes to the performance improvement of 1.4% mAP and 1.5% NDS, confirming the effectiveness of BEV-Paste.

4.3.4 Multi-Scale Cross-Task Head

The MSCT head contains the Multi-Task Distillation (MTD) module and the Dual Supervision (DS) on prediction from Multi-Scale (MS) image features. A number of experiments are carried out to further verify the effectiveness of each module and the results are shown in Table 6. Our MTD, DS and MS modules improve NDS performance by 0.4%, 0.5% and 1.0% respectively. We attribute this improvement to the fact that the multi-task distillation module supplements cross-task information, and the dual supervision further promotes the extraction and fusion of task-specific information and cross-task information, as well as the participation of multi-scale image features.

5. Conclusion and Discussion

In this paper, we propose SA-BEV to fully utilize the semantic information of images. SA-BEVPool filters out background virtual points and generates semantic-aware BEV features. BEV-Paste then pastes the semantic-aware BEV features of two frames to enhance data diversity. MSCT head introduces multi-task learning and facilitates the optimization of semantic-aware BEV features.

Our proposed components show strong universality. SA-BEVPool and BEV-Paste can be easily embedded into most BEV-based detectors and bring stable improvements. Besides, we believe that introducing multi-task learning into the generation of semantic-aware BEV features adds a valuable perspective and will inspire future works.

Still, there are limitations in SA-BEV. The thresholds used in SA-BEVPool are manually set, making it hard to achieve optimal performance. BEV-Paste may cause incorrect object overlaps and occlusions when pasting the semantic-aware BEV feature of one frame to another. Those are what we will tackle next. We also would like to extend SA-BEV into a multi-modal detector to activate the complementarity between image and LiDAR.

Acknowledgment

This paper was supported by National Natural Science Foundation of China under grants 62176017 and U20B2069.

References

- [1] 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 multi-modal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020.
- [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 213–229. Springer, 2020.
- [3] Jaeseok Choi, Yeji Song, and Nojun Kwak. Part-aware data augmentation for 3d object detection in point cloud. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3391–3397. IEEE, 2021.
- [4] 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 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 1201–1209, 2021.
- [5] Di Feng, Yiyang Zhou, Chenfeng Xu, Masayoshi Tomizuka, and Wei Zhan. A simple and efficient multi-task network for 3d object detection and road understanding. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7067–7074. IEEE, 2021.
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [7] Peiyun Hu, Jason Zizlar, David Held, and Deva Ramanan. What you see is what you get: Exploiting visibility for 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11001–11009, 2020.
- [8] Junjie Huang and Guan Huang. Bevdet4d: Exploit temporal cues in multi-camera 3d object detection. *arXiv preprint arXiv:2203.17054*, 2022.
- [9] Junjie Huang and Guan Huang. Bevpoolv2: A cutting-edge implementation of bevdet toward deployment. *arXiv preprint arXiv:2211.17111*, 2022.
- [10] Junjie Huang, Guan Huang, Zheng Zhu, Ye Yun, and Dalong Du. Bevdet: High-performance multi-camera 3d object detection in bird-eye-view. *arXiv preprint arXiv:2112.11790*, 2021.
- [11] Yanqin Jiang, Li Zhang, Zhenwei Miao, Xiatian Zhu, Jin Gao, Weiming Hu, and Yu-Gang Jiang. Polarformer: Multi-camera 3d object detection with polar transformers. *arXiv preprint arXiv:2206.15398*, 2022.
- [12] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12697–12705, 2019.
- [13] Youngwan Lee, Joong-Won Hwang, Sangrok Lee, Yuseok Bae, and Jongyoul Park. An energy and gpu-computation efficient backbone network for real-time object detection. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 752–760. Computer Vision Foundation / IEEE, 2019.
- [14] Yin hao Li, Han Bao, Zheng Ge, Jinrong Yang, Jianjian Sun, and Zeming Li. Bevstereo: Enhancing depth estimation in multi-view 3d object detection with dynamic temporal stereo. *arXiv preprint arXiv:2209.10248*, 2022.
- [15] Yin hao Li, Zheng Ge, Guanyi Yu, Jinrong Yang, Zengran Wang, Yukang Shi, Jianjian Sun, and Zeming Li. Bevdepth: Acquisition of reliable depth for multi-view 3d object detection. *arXiv preprint arXiv:2206.10092*, 2022.
- [16] Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng Dai. Bevformer: Learning bird’s-eye-view representation from multi-camera images via spatiotemporal transformers. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part IX*, pages 1–18. Springer, 2022.
- [17] Ming Liang, Bin Yang, Yun Chen, Rui Hu, and Raquel Urtasun. Multi-task multi-sensor fusion for 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [18] Yingfei Liu, Tiancai Wang, Xiangyu Zhang, and Jian Sun. Petr: Position embedding transformation for multi-view 3d object detection. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXVII*, pages 531–548. Springer, 2022.
- [19] Yingfei Liu, Junjie Yan, Fan Jia, Shuailin Li, Qi Gao, Tiancai Wang, Xiangyu Zhang, and Jian Sun. Petr v2: A unified framework for 3d perception from multi-camera images. *arXiv preprint arXiv:2206.01256*, 2022.
- [20] Zhijian Liu, Haotian Tang, Alexander Amini, Xingyu Yang, Huizi Mao, Daniela Rus, and Song Han. Bevfusion: Multi-task multi-sensor fusion with unified bird’s-eye view representation. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2023.
- [21] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*.
- [22] Jonah Philion and Sanja Fidler. Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d. In *Proceedings of the European Conference on Computer Vision*, 2020.
- [23] Zhile Ren and Erik B Sudderth. 3d object detection with latent support surfaces. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 937–946, 2018.
- [24] Khailash Santhakumar, B Ravi Kiran, Thomas Gauthier, Senthil Yogamani, et al. Exploring 2d data augmentation for 3d monocular object detection. *arXiv preprint arXiv:2104.10786*, 2021.
- [25] Simon Vandenhende, Stamatios Georgoulis, Wouter Van Gansbeke, Marc Proesmans, Dengxin Dai, and Luc Van Gool. Multi-task learning for dense prediction tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.

- [26] Simon Vandenhende, Stamatios Georgoulis, and Luc Van Gool. Mti-net: Multi-scale task interaction networks for multi-task learning. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16*, pages 527–543. Springer, 2020.
- [27] Sourabh Vora, Alex H Lang, Bassam Helou, and Oscar Beijbom. Pointpainting: Sequential fusion for 3d object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4604–4612, 2020.
- [28] Chunwei Wang, Chao Ma, Ming Zhu, and Xiaokang Yang. Pointaugmenting: Cross-modal augmentation for 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11794–11803, 2021.
- [29] Tai Wang, ZHU Xinge, Jiangmiao Pang, and Dahua Lin. Probabilistic and geometric depth: Detecting objects in perspective. In *Conference on Robot Learning*, pages 1475–1485. PMLR, 2022.
- [30] Tai Wang, Xinge Zhu, Jiangmiao Pang, and Dahua Lin. FCOS3D: Fully convolutional one-stage monocular 3d object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, 2021.
- [31] Yue Wang, Vitor Campagnolo Guizilini, Tianyuan Zhang, Yilun Wang, Hang Zhao, and Justin Solomon. Detr3d: 3d object detection from multi-view images via 3d-to-2d queries. In *Conference on Robot Learning*, pages 180–191. PMLR, 2022.
- [32] Dan Xu, Wanli Ouyang, Xiaogang Wang, and Nicu Sebe. Pad-net: Multi-tasks guided prediction-and-distillation network for simultaneous depth estimation and scene parsing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [33] Yan Yan, Yuxing Mao, and Bo Li. Second: Sparsely embedded convolutional detection. *Sensors*, 2018.
- [34] Zetong Yang, Yanan Sun, Shu Liu, and Jiaya Jia. 3dssd: Point-based 3d single stage object detector. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11040–11048, 2020.
- [35] Tianwei Yin, Xingyi Zhou, and Philipp Krahenbuhl. Center-based 3d object detection and tracking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11784–11793, 2021.
- [36] Tianwei Yin, Xingyi Zhou, and Philipp Krähenbühl. Multi-modal virtual point 3d detection. *Advances in Neural Information Processing Systems*, 34:16494–16507, 2021.
- [37] Wenwei Zhang, Zhe Wang, and Chen Change Loy. Exploring data augmentation for multi-modality 3d object detection. *arXiv preprint arXiv:2012.12741*, 2020.
- [38] Yanan Zhang, Jiaxin Chen, and Di Huang. Cat-det: Contrastively augmented transformer for multi-modal 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 908–917, 2022.
- [39] Yanan Zhang, Di Huang, and Yunhong Wang. Pc-rgnn: Point cloud completion and graph neural network for 3d object detection. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 3430–3437, 2021.
- [40] Zhenyu Zhang, Zhen Cui, Chunyan Xu, Yan Yan, Nicu Sebe, and Jian Yang. Pattern-affinitive propagation across depth, surface normal and semantic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4106–4115, 2019.
- [41] Wu Zheng, Weiliang Tang, Li Jiang, and Chi-Wing Fu. Sessd: Self-ensembling single-stage object detector from point cloud. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14494–14503, 2021.
- [42] Chao Zhou, Yanan Zhang, Jiaxin Chen, and Di Huang. Octree-based transformer for 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5166–5175, 2023.
- [43] Xingyi Zhou, Vladlen Koltun, and Philipp Krähenbühl. Probabilistic two-stage detection. In *arXiv preprint arXiv:2103.07461*, 2021.
- [44] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019.
- [45] Benjin Zhu, Zhengkai Jiang, Xiangxin Zhou, Zeming Li, and Gang Yu. Class-balanced grouping and sampling for point cloud 3d object detection. *arXiv preprint arXiv:1908.09492*, 2019.

A. More Implementation Details

A.1. Data Augmentation

We augment both images and BEV features following the operation applied in [10]. For images, they are first down-sampled to the desired resolution. Then they are processed by random scaling with a range of [0.94, 1.11], random rotating with a range of $[-5.4^\circ, 5.4^\circ]$ and random flipping with a probability of 0.5. After that, the images are padded and cropped to a uniform shape. For BEV features, augmentation is applied on the virtual points whose features are cumulated to form BEV features. The coordinates of virtual points are processed by random scaling with a range of [0.95, 1.05], random flipping of the X and Y axes with a probability of 0.5 and random rotating with a range of $[-22.5^\circ, 22.5^\circ]$. Augmenting virtual points rather than BEV features themselves can generate more accurate augmented BEV features because the bilinear sampling is not required by the former. The additional BEV data augmentation (BDA) used by BEV-Paste also follows the above settings.

A.2. Detection Configuration

We use the detection head of CenterPoint [35] to detect 3D objects from semantic-aware BEV features and follow the settings used in BEVDepth [15]. The LiDAR coordinate system of nuScenes is used to represent the coordinate of points in the BEV space. The X and Y coordinates are in the range of $[-51.2, 51.2]$, and the Z coordinate is in the range of $[-5, 3]$. The BEV space is divided into pillars for cumulating virtual point features. When the resolution of input images is 256×704 , the pillars are in the size of $[0.8, 0.8, 8]$ and the BEV features are in the shape of 128×128 . For larger input images, the pillars are in the size of $[0.4, 0.4, 8]$ and the BEV features are in the shape of 256×256 .

B. More Experiment Results

We change the image backbone of SA-BEV to ResNet-101 when processing 512×1408 resolution images and compare it with other methods that also utilize ResNet-101 as their backbone. The results are shown in Table 7. SA-BEV achieves the best mAP and NDS, 2.9% and 1.4% higher than its baseline (i.e. BEVDepth [15]). It also exceeds other start-of-the-art methods that take 900×1600 resolution images as input. This comparison further proves the effectiveness of SA-BEV.

We also compare the detection precision of BEVDepth and SA-BEV in each category and show the results in Fig. 7. SA-BEV achieves better precision than BEVDepth in most of the categories. For instance, the APs on pedestrian and traffic cone are increased by about 10%, and the APs on car, truck, bus and bicycle are increased by about 3%. The

Table 7: Comparison with previous state-of-the-art multi-view 3D detectors on the nuScenes *val* set.

| Method | Backbone | Resolution | mAP \uparrow | NDS \uparrow |
|------------------|------------|-------------------|----------------|----------------|
| FCOS3D [30] | ResNet-101 | 900 \times 1600 | 0.343 | 0.415 |
| DETR3D [31] | ResNet-101 | 900 \times 1600 | 0.303 | 0.374 |
| PGD [29] | ResNet-101 | 900 \times 1600 | 0.369 | 0.428 |
| PETR [18] | ResNet-101 | 512 \times 1408 | 0.357 | 0.421 |
| BEVFormer [16] | ResNet-101 | 900 \times 1600 | 0.416 | 0.517 |
| PETrv2 [19] | ResNet-101 | 900 \times 1600 | 0.421 | 0.524 |
| PolarFormer [11] | ResNet-101 | 900 \times 1600 | 0.432 | 0.528 |
| BEVDepth [15] | ResNet-101 | 512 \times 1408 | 0.412 | 0.535 |
| SA-BEV | ResNet-101 | 512 \times 1408 | 0.441 | 0.549 |

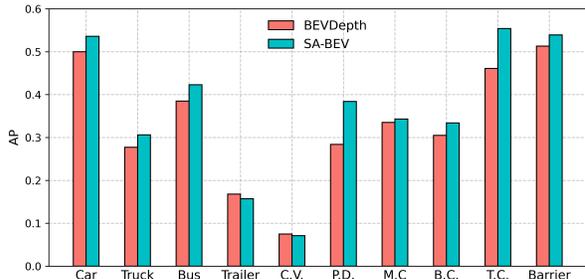


Figure 7: Comparison of BEVDepth and SA-BEV on AP for each category. C.V., P.D., M.C., B.C. and T.C. are the abbreviations of construction vehicle, pedestrian, motorcycle, bicycle and traffic cone respectively.

greater improvement in pedestrian and traffic cone categories indicates that the semantic-aware BEV features effectively preserve the information of small scale objects that is more likely to be submerged by the large proportion of background information.

C. More Visualization Results

We provide more visualization results of BEVDepth and SA-BEV in Fig. 8. With the help of semantic-aware BEV features, SA-BEV can recall objects in the far distance and identify the false truth precisely. Besides, SA-BEV generally predicts more accurate locations and directions of the objects, which is also important in actual practice.

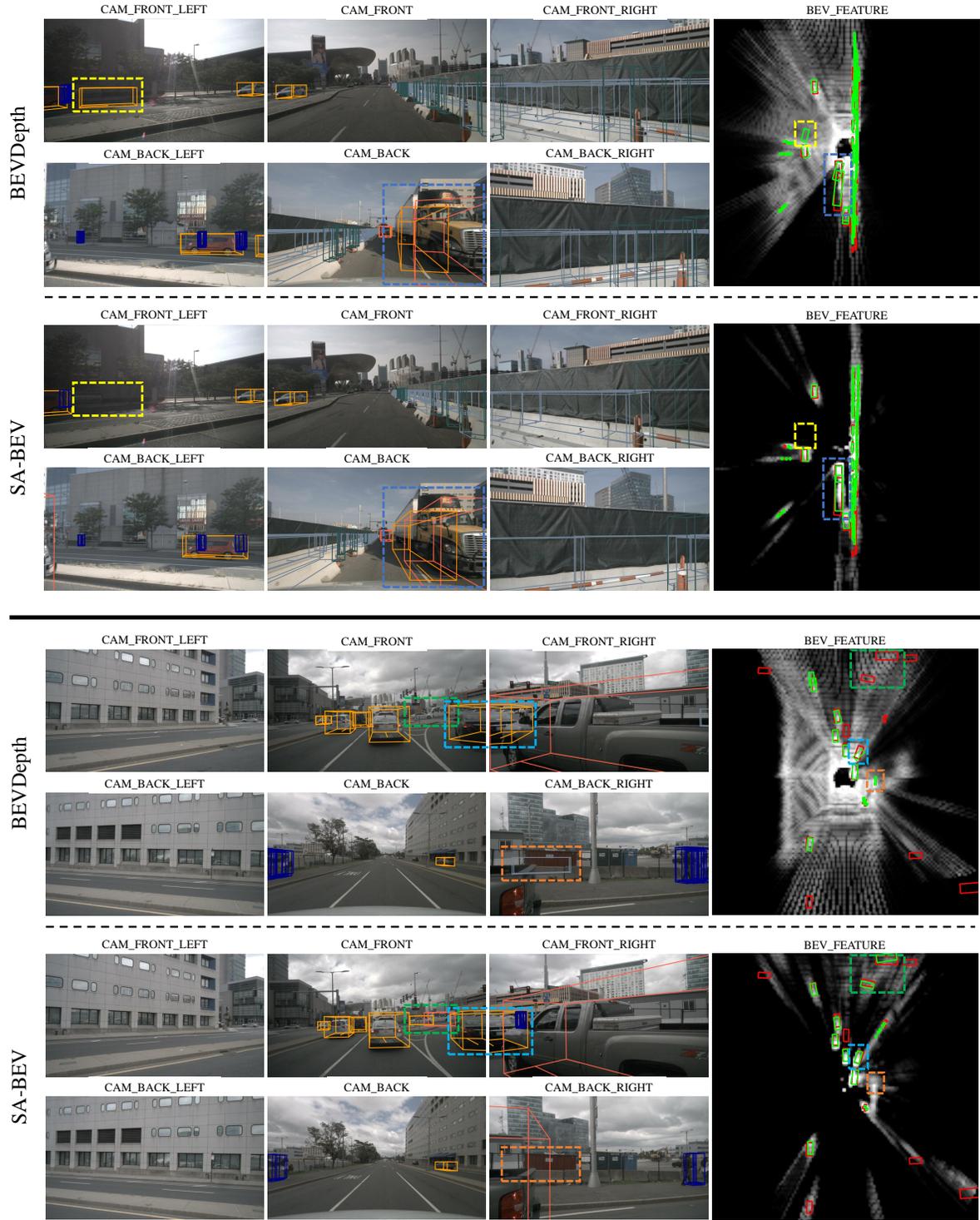


Figure 8: Visualization results on images and BEV features. The red boxes and green boxes on BEV features represent the ground truth and the predicted boxes, respectively. The dashed rectangles illustrate that the prediction of SA-BEV is more precise than BEVDepth.