

# UniVision: A Unified Framework for Vision-Centric 3D Perception

Yu Hong<sup>1</sup>    Qian Liu<sup>2</sup>    Huayuan Cheng<sup>1</sup>    Danjiao Ma<sup>2</sup>  
 Hang Dai<sup>3</sup>    Yu Wang<sup>2</sup>    Guangzhi Cao<sup>2</sup>    Yong Ding<sup>1</sup>

<sup>1</sup>Zhejiang University    <sup>2</sup>Pegasus Tech    <sup>3</sup>University of Glasgow

## Abstract

*The past few years have witnessed the rapid development of vision-centric 3D perception in autonomous driving. Although the 3D perception models share many structural and conceptual similarities, there still exist gaps in their feature representations, data formats, and objectives, posing challenges for unified and efficient 3D perception framework design. In this paper, we present UniVision, a simple and efficient framework that unifies two major tasks in vision-centric 3D perception, i.e., occupancy prediction and object detection. Specifically, we propose an explicit-implicit view transform module for complementary 2D-3D feature transformation. We propose a local-global feature extraction and fusion module for efficient and adaptive voxel and BEV feature extraction, enhancement, and interaction. Further, we propose a joint occupancy-detection data augmentation strategy and a progressive loss weight adjustment strategy which enables the efficiency and stability of the multi-task framework training. We conduct extensive experiments for different perception tasks on four public benchmarks, including nuScenes LiDAR segmentation, nuScenes detection, OpenOccupancy, and Occ3D. UniVision achieves state-of-the-art results with +1.5 mIoU, +1.8 NDS, +1.5 mIoU, and +1.8 mIoU gains on each benchmark, respectively. We believe that the UniVision framework can serve as a high-performance baseline for the unified vision-centric 3D perception task. The code will be available at <https://github.com/Cc-Hy/UniVision>.*

## 1. Introduction

3D perception is the primary task in autonomous driving systems, and its purpose is to use the data obtained by a series of sensors (e.g., LiDAR, Radar, and camera) to derive a comprehensive understanding of the driving scenes, which is used for the subsequent planning and decision-making. In the past, the field of 3D perception has been dominated

by LiDAR-based models due to the accurate 3D information from point cloud data. However, LiDAR-based systems are costly, vulnerable to bad weather, and inconvenient to deploy. In comparison, vision-based systems have many advantages such as low cost, easy deployment, and good scalability. Thus, vision-centric 3D perception has attracted extensive attention from the researchers.

Recently, vision-based 3D detection has been significantly improved via feature representation transformation [20, 28, 40], temporal fusion [14, 47, 67], and supervision signal design [8, 9, 18], continuously narrowing the gap with LiDAR-based models. Lately, vision-based occupancy task has witnessed rapid development [3, 21, 43, 51, 53, 62]. Unlike using 3D bounding boxes to represent some whitelist objects, occupancy can more comprehensively describe the geometric and semantics of the driving scene and it is less limited to the shape and category of objects.

Although the detection methods and the occupancy methods share many structural and conceptual similarities, it is not well investigated to simultaneously tackle the two tasks and explore the interrelationship between them. Occupancy models and detection models often extract different feature representations. The occupancy prediction task requires exhaustive semantic and geometric judgments across different spatial positions, so the voxel representation is widely used to preserve fine-grained 3D information. In the detection task, the BEV representation is preferred since most objects are on the same horizontal level with minor overlap. Compared to the BEV representation, the voxel representation is elaborate but less efficient. Also, many advanced operators (e.g., shifted window attention [34] and deformable convolution [10]) are primarily designed and optimized for 2D features, making their integration with the 3D voxel representation less straightforward. The BEV representation is more time-efficient and memory-efficient but it is sub-optimal for dense spatial prediction as it loses structural information in the height dimension. Apart from feature representations, different perception tasks also differ in their data formats and objectives. Thus, it is a great

challenge to ensure the unity and efficiency of training a multi-task 3D perception framework.

To further exploit the potential of vision models and explore the correlations between different perception tasks, we present UniVision, a unified framework that simultaneously handles the tasks of 3D detection and occupancy prediction. The framework follows a join-divide-join diagram. Given the surrounding images as input, we use a shared network for image feature extraction. We propose a novel view transformation module that combines both depth-guided lifting and query-guided sampling for complementary 2D-3D feature transformation. After that, the network splits into voxel-based and BEV-based branches in parallel to extract features with local and global receptive field awareness, leveraging the advantages of different feature representations. We then employ adaptive feature interaction between the two feature representations to enhance each other, followed by task-specific heads for different perception tasks. In addition to the framework design, we present a joint occupancy-detection data augmentation method and a multi-task training strategy for the efficient training of the UniVision framework. We conduct extensive experiments on four benchmarks, including nuScenes LiDAR segmentation [2], nuScenes detection [2], OpenOccupancy [51], and Occ3D [43]. The proposed UniVision framework not only efficiently handles different 3D perception tasks, but also achieves state-of-the-art performance on different benchmarks.

Our contributions can be summarized as: **i)** We propose a simple and efficient framework for unified vision-centric 3D perception, which simultaneously handles the detection and occupancy tasks. Extensive experiments demonstrate the generalization and superiority of the UniVision framework with state-of-the-art performance on different benchmarks. **ii)** We propose an explicit-implicit (Ex-Im) view transform method that combines both depth-guided lifting and query-guided sampling, facilitating complementary 2D-3D feature transformation. **iii)** We propose a local-global feature extraction and fusion module for efficient and adaptive feature extraction, enhancement, and interaction. **iv)** We present a joint Occupancy-Detection (Occ-Det) data augmentation method and a progressive loss weight adjustment strategy to enable efficient training of the multi-task framework.

## 2. Related Works

### 2.1. Vision-based 3D Detection

Vision-based 3D detection aims to locate and classify 3D objects with images from single or multiple cameras. Early methods [1, 6, 36, 38, 49, 60] extend advanced 2D object detection methods [11, 44] to the 3D case by predicting additional 3D attributes based on 2D ones. Later, Bird’s-Eye-View (BEV) based diagram has become the mainstream.

CaDDN [40] leverages the Lift-Splat-Shoot (LSS) [39] diagram to transform monocular images into BEV features, executing the detection process within BEV frameworks [25, 56, 58]. BEVDet [20] and BEVFormer [28] transform images from surround-view cameras into a single BEV feature map for full-range detection. Besides, various approaches [32, 52, 59] make efforts to introduce the DETR diagram into the 3D area. Recent methods have further improved the performance of 3D object detection from the perspective of long-term temporal fusion [14, 47, 67] and sparse representations [30, 31].

### 2.2. Vision-based Occupancy Prediction

Occupancy prediction, also known as semantic scene completion (SSC), requires exhaust judgments for positions in the scene, including whether the position is occupied and the category of occupation. Early vision-based methods [4, 13, 42, 54] use images enriched with additional geometric information, such as RGB-D images, to execute occupancy prediction. MonoScene [3] is the first method to infer dense geometry and semantics from a single monocular image. TPVFormer [21] enhances the commonly utilized Bird’s-Eye-View (BEV) by introducing the Tri-Perspective-View (TPV), thereby augmenting the representation with Z-axis information. OccFormer [62] proposes a dual-path transformer network to process the 3D volume for semantic occupancy prediction. Also, works like OpenOccupancy [51], Occ3D [43] and SurroundOcc [53] propose pipelines for generating high-quality dense occupancy labels.

### 2.3. Multi-task Framework

Multi-task frameworks [5, 23, 37, 57, 63] strive to efficiently manage various tasks within a singular network. In the 2D vision area, Mask-RCNN [16] proposes a unified network for object detection and mask segmentation. UberNet [24] simultaneously handles a variety of low, medium, and high-level visual tasks in an end-to-end manner. In LiDAR-based 3D perception, initiatives such as LidarMTL [12] and LidarMultiNet [57] leverage a shared network for tasks encompassing 3D detection, segmentation, and road understanding. A major advantage of multi-task networks is to save computational and storage overheads by utilizing shared model structure and weights. However, the performance of individual tasks frequently diminishes as the network navigates trade-offs between different objectives, posing challenges for multi-task framework design.

## 3. Method

### 3.1. Framework Overview

Fig. 1 shows the overall architecture of the UniVision framework. Given multi-view images  $\{I^i | I^i \in \mathbb{R}^{W_I \times H_I \times 3}\}$ ,  $i \in [1, N]$  from the surrounding  $N$  cameras

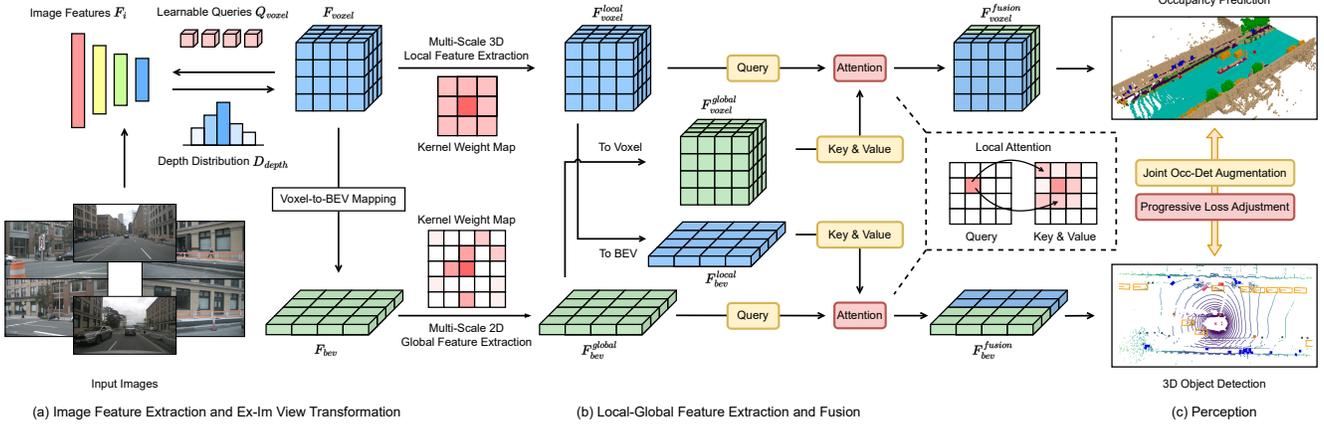


Figure 1. The overall architecture of UniVision. After extracting image features from the inputs, we use the Ex-Im view transform module for complementary 2D-3D feature transformation. We then propose the local-global feature extraction and fusion block for adaptive BEV and voxel feature extraction, enhancement, and interaction, which are attached to task-specific perception heads. During training, the joint Occ-Det augmentation and progressive loss weight adjustment strategy are equipped for efficient multi-task training.

as input, an image feature extraction network is first utilized to extract image features  $F_{img}$  from them. The 2D image features  $F_{img}$  are then lifted to 3D voxel features  $F_{voxel}$  with the Ex-Im view transform module, which combines depth-guided explicit feature lifting and query-guided implicit feature sampling. The voxel features  $F_{voxel}$  are sent into the local-global feature extraction and fusion block to extract local-context-aware voxel features and global-context-aware BEV features respectively. Then, we employ the cross-representation feature interaction module to perform information exchange on the voxel features and BEV features, which are used for different downstream perception tasks. During training, the joint Occ-Det data augmentation and the progressive loss weight adjustment strategy are used for efficient training of the UniVision framework.

### 3.2. Ex-Im View Transform

**Depth-guided Explicit Feature Lifting.** Following the Lift-Splat-Shoot (LSS) [39] diagram, we perform the voxel pooling operation [20] on the per-pixel depth distribution  $D_{depth} \in \mathbb{R}^{D \times H \times W}$  and the image features  $F_{img} \in \mathbb{R}^{C \times H \times W}$  to extract the voxel features:

$$F_{voxel}^{ex} = VoxelPooling(D_{depth}, F_{img}) \quad (1)$$

Since  $F_{voxel}^{ex}$  is generated with the explicit depth distribution estimation, we refer to it as the explicit voxel features.

**Query-guided Implicit Feature Sampling.** However,  $F_{voxel}^{ex}$  has some defects in representing the 3D information. The accuracy of  $F_{voxel}^{ex}$  is highly related to the accuracy of the estimated depth distribution  $D_{depth}$ . Also, the points generated from LSS are unevenly distributed. Points are dense close to the camera and are sparse in distance. We thus further use the query-guided feature sampling to compensate for the above shortcomings of  $F_{voxel}^{ex}$ . We define

the learnable voxel queries  $q_{voxel} \in \mathbb{R}^{C \times X \times Y \times Z}$ , and use a 3D transformer to sample the features from the images. For each voxel query, we project its center  $\mathbf{c}$  onto the image plane with calibration matrix  $\mathbf{P}$  for the reference point  $\mathbf{p}$ , and then use  $N$  transformer blocks. Each block includes a deformable cross-attention (DCA) [66] layer, a 3D convolution (Conv) layer, and a feed-forward network (FFN) [45]:

$$\mathbf{p} = \mathbf{P} \times \mathbf{c} \quad (2)$$

$$q^{i+1} = FFN(Conv(DCA(q^i, \mathbf{p}, F_{img}))) \quad (3)$$

$$F_{voxel}^{im} = q^N \quad (4)$$

Compared to the points generated from LSS, the voxel queries are evenly distributed in the 3D space and they are learned from the statistical properties of all training samples, which is independent of the depth prior information used in LSS. Thus,  $F_{voxel}^{ex}$  and  $F_{voxel}^{im}$  complement each other, and we concatenate them as the output features of the view transform module:

$$F_{voxel} = F_{voxel}^{ex} \parallel F_{voxel}^{im} \quad (5)$$

where  $\parallel$  is the concatenate operation. The Ex-Im view transform module is illustrated in Fig. 2.

### 3.3. Local-Global Feature Extraction and Fusion

Given input voxel features  $F_{voxel} \in \mathbb{R}^{C \times X \times Y \times Z}$ , we first stack the features in  $Z$  axis and use a convolution layer to reduce the channels to obtain the BEV features  $F_{bev} \in \mathbb{R}^{C \times X \times Y}$ :

$$F_{bev} = Conv(Stack(F_{voxel}, dim = Z)) \quad (6)$$

Then, the model splits into two parallel branches for feature extraction and enhancement.

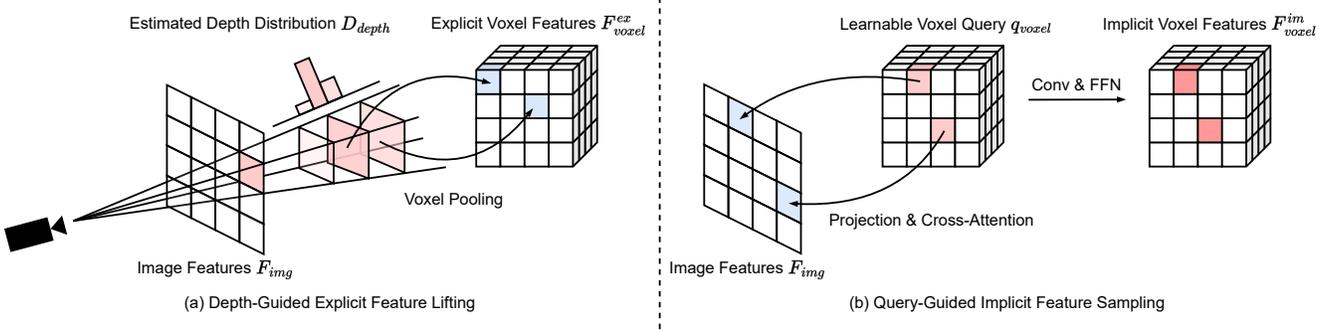


Figure 2. The Ex-Im view transform module. (a) Depth-guided Explicit Feature Lifting. (b) Query-guided Implicit Feature Sampling.

**Local feature extraction.** For  $F_{voxel}$ , we use a local feature extraction branch composed of 3D convolutions to extract local features of each spatial position. We extend ResNet [17] to ResNet3D to extract multi-scale voxel features  $\{F_{voxel}^i | F_{voxel}^i \in \mathbb{R}^{C \cdot 2^i \times \frac{X}{2^i} \times \frac{Y}{2^i} \times \frac{Z}{2^i}}\}$  from  $F_{voxel}$ . We then use the FPN [29] structure from SECOND [56] in 3D version to merge  $\{F_{voxel}^i\}$  into  $F_{voxel}^{local} \in \mathbb{R}^{C \times X \times Y \times Z}$ .

**Global feature extraction.** The BEV features  $F_{bev}$  retain the information at the object level and they are computationally efficient. Thus, We propose a global feature extraction branch to extract features with a global receptive field based on the BEV representation. We use a network with deformable convolution v3 (DCNV3) [50] to dynamically gather global information for multi-scale BEV features  $\{F_{bev}^i | F_{bev}^i \in \mathbb{R}^{C \cdot 2^i \times \frac{X}{2^i} \times \frac{Y}{2^i}}\}$ . And  $\{F_{bev}^i\}$  goes through the SECOND FPN structure for the merged BEV feature  $F_{bev}^{global} \in \mathbb{R}^{C \times X \times Y}$ .

**Cross-Representation Feature Interaction.** After generating the local-context-aware voxel features  $F_{voxel}^{local}$  and the global-context-aware BEV features  $F_{bev}^{global}$  from the input voxel features, we use the cross-representation feature interaction module to enable adaptive information exchange between the two feature representations for further enhancement. We first map the BEV features  $F_{bev}^{global}$  to voxel features  $F_{voxel}^{global}$  and map voxel features  $F_{voxel}^{local}$  to BEV features  $F_{bev}^{local}$  by Z-dimensional repetition or addition:

$$F_{voxel}^{global} = \text{repeat}(F_{bev}^{global}, \text{dim} = Z) \quad (7)$$

$$F_{bev}^{local} = \text{add}(F_{voxel}^{local}, \text{dim} = Z) \quad (8)$$

For the voxel representation, we use  $F_{voxel}^{local}$  from the voxel branch as the query, and  $F_{voxel}^{global}$  from the BEV branch as the key and value. And we generalize the neighborhood attention transformer [15] from self-attention to cross-attention to perform information gathering within a local perception field  $\Delta p$ , and a symmetric process is applied on the BEV

features:

$$F_{voxel}^{fusion} = \text{Attn}(q = F_{voxel}^{local}, k \& v = F_{voxel}^{global}, \Delta p) \quad (9)$$

$$F_{bev}^{fusion} = \text{Attn}(q = F_{bev}^{global}, k \& v = F_{bev}^{local}, \Delta p) \quad (10)$$

Specifically, we set  $\Delta p$  to  $3 \times 3$  for the BEV features and  $3 \times 3 \times 3$  for the voxel features.

### 3.4. Heads and Losses

We attach task-specific heads to  $F_{voxel}^{fusion}$  and  $F_{bev}^{fusion}$  for different perception tasks. For the occupancy task, we use two fully connected layers to map the feature channels to the number of occupancy categories. We follow the loss setting in OpenOccupancy [51], which combines the cross-entropy loss, Lovasz softmax loss, geometry affinity loss, and semantic affinity loss:

$$\mathcal{L}_{occ} = \lambda_1 \cdot L_{ce} + \lambda_2 \cdot L_{lovasz} + \lambda_3 \cdot L_{geo} + \lambda_4 \cdot L_{sem} \quad (11)$$

For the detection task, we use the center-based head [58] and the detection loss is composed of the classification loss  $L_{cls}$  and the regression loss  $L_{reg}$ :

$$\mathcal{L}_{det} = \lambda_5 \cdot L_{cls} + \lambda_6 \cdot L_{reg} \quad (12)$$

Also, we add the depth loss  $L_{depth}$  used in BEVDepth [27] as the image level supervision:

$$\mathcal{L}_{img} = \lambda_7 \cdot L_{depth} \quad (13)$$

The overall loss function is defined as:

$$\mathcal{L} = \mathcal{L}_{img} + \mathcal{L}_{det} + \mathcal{L}_{occ} \quad (14)$$

**Progressive Loss Weight Adjustment Strategy.** In practice, we find that directly combining the above losses tends to lead to a failed training process and the network cannot converge. In the early stage of training, the voxel features  $F_{voxel}$  are randomly distributed, and the supervision in the occupancy head and detection head contribute less than other losses in the convergence. Meanwhile, the

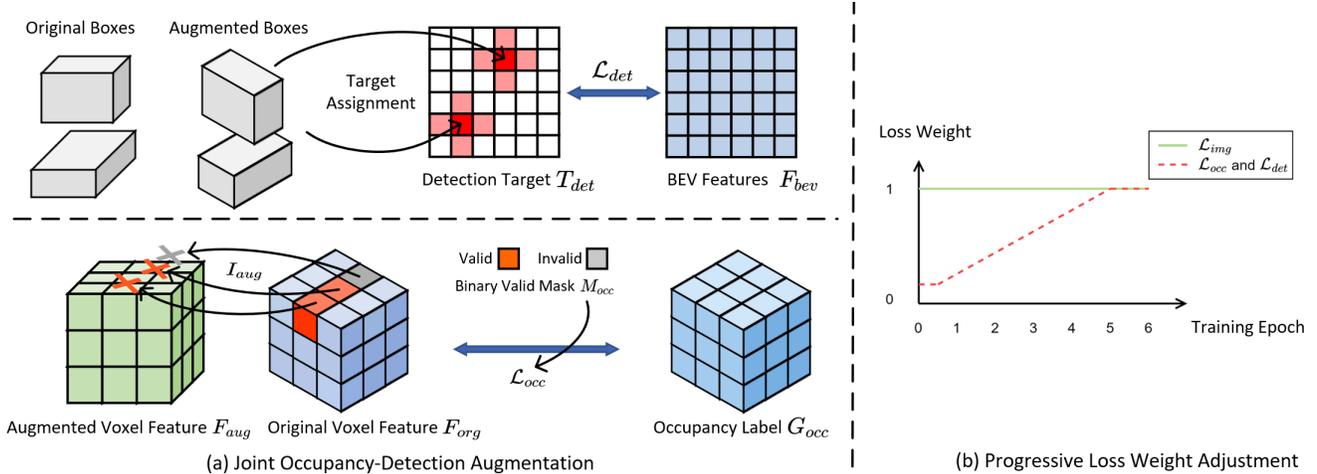


Figure 3. Illustration of (a) joint occupancy-detection augmentation and (b) progressive loss weight adjustment strategy.

loss items like the classification loss  $\mathcal{L}_{cls}$  in the detection task are very large and dominate the training process, making the model difficult to optimize.

To overcome this, we propose the progressive loss weight adjustment strategy to dynamically adjust the loss weights. Specifically, a control parameter  $\delta$  is added to the non-image-level losses, *i.e.*, the occupancy loss and the detection loss, to adjust loss weights in different training periods. The control weight  $\delta$  is set to a small value  $V_{min}$  at the beginning and gradually increase to  $V_{max}$  in  $N$  training epochs:

$$\mathcal{L} = \mathcal{L}_{img} + \delta \cdot \mathcal{L}_{det} + \delta \cdot \mathcal{L}_{occ} \quad (15)$$

$$\delta = \max(V_{min}, \min(V_{max}, \frac{i}{N} \cdot V_{max})) \quad (16)$$

where  $i$  denotes the  $i_{th}$  training epoch. In this case, the optimization process first focuses on the image-level information (depth and semantics) to generate reasonable voxel representations, and then on the subsequent perception tasks. The progression is illustrated in Fig. 3 (b).

### 3.5. Joint Occ-Det Spatial Data Augmentation

In the 3D detection task, spatial-level data augmentation is also effective in improving model performance in addition to the common image-level data augmentation. However, it is not straightforward to apply spatial-level augmentation in the occupancy task. When we apply data augmentation such as random scaling and rotation to the discrete occupancy labels  $G_{occ} \in \mathbb{R}^{X \times Y \times Z}$ , it is difficult to determine the semantics of the generated voxels. Thus, the existing methods only apply simple spatial augmentation like random flipping in occupancy tasks.

To address this problem, we propose a joint Occ-Det spatial data augmentation to allow simultaneous augmentation

in both the 3D detection task and the occupancy task in our UniVision framework. Since the 3D box labels are in continuous values and the augmented 3D box can be directly calculated for training, we follow the augmentation method for detection from BEVDet [20]. Although the occupancy labels are discrete and difficult to operate on, the voxel features can be regarded as continuous and can be handled with operations like sampling and interpolation. Thus, we propose to transform the voxel features instead of directly operating on the occupancy labels for the data augmentation.

Specifically, we first sample spatial data augmentations and calculate the corresponding 3D transformation matrix  $\mathbf{M}_{aug}$ . For the occupancy labels  $G_{occ} \in \mathbb{R}^{X \times Y \times Z}$  and its voxel indices  $\mathbf{I}_{org} \in \mathbb{R}^{X \times Y \times Z \times 3}$ , we compute the their 3D coordinates  $\mathbf{C}_{org} \in \mathbb{R}^{X \times Y \times Z \times 3}$ . We then apply  $\mathbf{M}_{aug}$  to  $\mathbf{C}_{org}$ , and normalize them to obtain the voxel indices  $\mathbf{I}_{aug}$  in the augmented voxel features:

$$\mathbf{C}_{org} = \mathbf{P}_{i-c} \times \mathbf{I}_{org} \quad (17)$$

$$\mathbf{I}_{aug} = \mathbf{P}_{c-i} \times \mathbf{M}_{aug} \times \mathbf{C}_{org} \quad (18)$$

where  $\mathbf{P}_{i-c}$  and  $\mathbf{P}_{c-i}$  are the transformation matrices between voxel indices and spatial coordinates. Then, we sample the voxel features  $F_{aug}$  with the voxel indices  $\mathbf{I}_{aug}$ :

$$F_{org} = S(F_{aug}, \mathbf{I}_{org}) \quad (19)$$

where  $S$  denotes sampling operation and  $F_{org}$  are the sampled voxel features that correspond to the original occupancy labels  $G_{occ}$  without transformation, which can be used for loss calculation. Noticeably, some sampling positions can fall out of range, and we ignore these voxels by adding a binary mask  $M_{occ} \in \{0, 1\}^{X \times Y \times Z}$  when calculating the occupancy losses:

$$\mathcal{L}_{occ} = f(G_{occ}, F_{org}) \times M_{occ} \quad (20)$$

where  $f$  denotes the loss functions. We illustrate the joint Occ-Det augmentation in Fig. 3 (a).

## 4. Experiments and Discussions

### 4.1. Datasets and Evaluation Metrics

**NuScenes.** NuScenes [2] is a modern and large-scale dataset for autonomous driving which contains 1000 driving scenes. It provides sensor data including LiDAR, radar, camera, and support benchmarking for different autonomous driving tasks, *e.g.*, 3D detection, LiDAR segmentation, and motion planning.

**NuScenes LiDAR Segmentation.** Following recent OccFormer [62] and TPVFormer [21], we use camera images as input for the LiDAR segmentation task, and the LiDAR data is only used to provide 3D positions for querying the output features. We use the mean intersection over union (mIoU) as the evaluation metric.

**NuScenes 3D Object Detection.** For the detection task, we use the official metrics of nuScenes, the nuScenes Detection Score (NDS), which is a weighted sum of mean Average Precision (mAP) and several metrics, including Average Translation Error (ATE), Average Scale Error (ASE), Average Orientation Error (AOE), Average Velocity Error (AVE) and Average Attribute Error (AAE).

**OpenOccupancy.** The OpenOccupancy benchmark [51] is based on the nuScenes dataset and provides semantic occupancy labels of  $512 \times 512 \times 40$  resolution. The labeled classes are the same as those in the LiDAR segmentation task and we use mIoU as the evaluation metric.

**Occ3D.** The Occ3D benchmark [43] is based on the nuScenes dataset and provides semantic occupancy labels of  $200 \times 200 \times 16$  resolution. Occ3D further provides visible masks for training and evaluation. The labeled classes are the same as those in the LiDAR segmentation task and we use mIoU as the evaluation metric.

### 4.2. Experimental Settings

**NuScenes LiDAR Segmentation** For the LiDAR segmentation task, we use the sparse LiDAR segmentation labels as the supervision only and no extra dense occupancy labels from other benchmarks are used. We use ResNet-101 [17] as the image backbone and the image resolution is set to  $896 \times 1600$ . The model is trained for 20 epochs with a total batch size of 32. We use the AdamW [35] optimizer and the learning rate is set to 0.0002. No temporal information or test time augmentation (TTA) is used.

**NuScenes Detection** For the results on the nuScenes detection benchmark [2], we provide two versions of comparison results. In the first version, we select three previous best methods [20,27,28] and align the training settings including image backbone, input resolution, batch size, and learning rate to make a fair comparison. We use the ResNet-50 im-

age backbone and the models are trained for 20 epochs with a learning rate of 0.0002. The batch size is set to 32 when the input resolution is  $256 \times 704$  and 16 when the input resolution is  $512 \times 1408$ .

In the second version, we upscale the model and compare the results with those reported in other methods' papers. We use ResNet-101 as the image backbone and the image resolution is set to  $512 \times 1408$ . The model is trained for 20 epochs with a total batch size of 32 using the AdamW optimizer and the learning rate is set to 0.0002. For UniVision4D, we initialize the model weights from the single-frame version and train the model for 10 epochs with a learning rate of 0.0001. Notably, we do not use the CBGS [65] strategy in other methods.

**OpenOccupancy** We use the ResNet-50 image backbone and the models are trained for 20 epochs with a learning rate of 0.0002, and the batch size is set to 32. The input resolution is set to  $512 \times 1408$ . Considering that the OpenOccupancy benchmark provides labels with a resolution of  $512 \times 512 \times 40$  which is very memory-consuming, we down-sample the labels to half the original resolution for training. During the inference phase, we up-sample the output to the original resolution. No temporal information or TTA is used.

**Occ3D** Considering that the Occ3D benchmark [43] is newly released with few reported results, we use the official codes of the compared methods [20,27,28,51,53] and align the training settings including image backbone, input resolution, batch size, and learning rate for a fair comparison. We choose the ResNet-50 backbone and train the models for 20 epochs. The batch size is set to 32 when the input resolution is  $256 \times 704$  and 16 when the input resolution is  $512 \times 1408$ . We use the AdamW optimizer and the learning rate is set to 0.0002. During training and inference, the camera visible mask is used for loss calculation or metric evaluation. No temporal information or TTA is used.

**Ablation Studies** For the ablation studies, we use the nuScenes detection benchmark [2] and the Occ3D benchmark [43] to validate the effectiveness of the components. We use the ResNet-50 backbone and the image resolution is set to  $256 \times 704$ . All the models are trained for 20 epochs with a total batch size of 32. We use the AdamW optimizer and the learning rate is set to 0.0002.

### 4.3. Results

**NuScenes LiDAR Segmentation.** We show the results on nuScenes LiDAR Segmentation benchmark in Tab. 1. UniVision significantly surpasses state-of-the-art (SOTA) vision-based method OccFormer [62] by 1.5 mIoU and sets a new record among vision-based models on the leaderboard. Notably, UniVision also outperforms some LiDAR-based models, such as PolarNet [61] and DB-UNet [46].

**NuScenes 3D Object Detection.** As shown in Tab. 2,

Method	Modality	mIoU	Class															
			barrier	bycicle	bus	car	cons. vehi.	motorcycle	pedestrian	traffic_cone	trailer	truck	driv. surf.	other flat	sidewalk	terrain	manmade	vegetation
PolarNet [61]	LiDAR	69.4	72.2	16.8	77.0	86.5	51.1	69.7	64.8	54.1	69.7	63.5	96.6	67.1	77.7	72.1	87.1	84.5
DB-UNet [46]	LiDAR	70.1	67.5	23.7	75.3	82.1	47.0	72.5	67.3	66.6	74.3	60.1	96.9	64.2	76.9	73.4	88.0	86.7
PolarStream [61]	LiDAR	73.4	71.4	27.8	78.1	82.0	61.3	77.8	75.1	72.4	79.6	63.7	96.0	66.5	76.9	73.0	88.5	84.8
Cylinder3D++ [64]	LiDAR	77.9	82.8	33.9	84.3	89.4	69.6	79.4	77.3	73.4	84.6	69.4	97.7	70.2	80.3	75.5	90.4	87.6
LidarMultiNet [57]	LiDAR	81.4	80.4	48.4	94.3	90.0	71.5	87.2	85.2	80.4	86.9	74.8	97.8	67.3	80.7	76.5	92.1	89.6
TPVFormer [21]	Camera	69.4	<b>74.0</b>	27.5	86.3	85.5	<b>60.7</b>	68.0	62.1	49.1	81.9	68.4	94.1	59.5	66.5	63.5	83.8	79.9
OccFormer [62]	Camera	70.8	72.8	29.9	<b>87.9</b>	85.6	57.1	74.9	63.2	53.5	83.0	67.6	94.8	61.9	70.0	<b>66.0</b>	<b>84.0</b>	<b>80.5</b>
UniVision	Camera	<b>72.3</b>	72.1	<b>34.0</b>	85.5	<b>89.5</b>	59.3	<b>75.5</b>	<b>69.3</b>	<b>65.8</b>	<b>84.2</b>	<b>71.4</b>	<b>96.1</b>	<b>67.4</b>	<b>71.9</b>	65.0	77.9	71.7

Table 1. Results on nuScenes LiDAR segmentation benchmark (test). UniVision sets a new record on the leaderboard among camera-based methods and has comparable or better results than some of the LiDAR-based methods.

Method	Backbone	Resolution	mAP $\uparrow$	NDS $\uparrow$	mATE $\downarrow$	mASE $\downarrow$	mAOE $\downarrow$	mAVE $\downarrow$	mAAE $\downarrow$
BEVFormer* [28]	R50	256 $\times$ 704	0.232	0.339	0.852	0.308	0.720	<b>0.648</b>	0.240
BEVDet* [20]	R50	256 $\times$ 704	0.278	0.348	0.783	0.289	0.686	0.860	0.289
BEVDepth* [27]	R50	256 $\times$ 704	0.292	0.386	0.724	<b>0.269</b>	<b>0.575</b>	0.796	0.241
UniVision*	R50	256 $\times$ 704	<b>0.312</b>	<b>0.397</b>	<b>0.682</b>	0.287	0.632	0.763	<b>0.224</b>
BEVFormer* [28]	R50	512 $\times$ 1408	0.292	0.388	0.819	0.294	0.635	<b>0.605</b>	<b>0.222</b>
BEVDet* [20]	R50	512 $\times$ 1408	0.329	0.394	0.724	0.275	0.607	0.852	0.247
BEVDepth* [27]	R50	512 $\times$ 1408	0.354	0.428	0.674	<b>0.267</b>	<b>0.506</b>	0.806	0.236
UniVision*	R50	512 $\times$ 1408	<b>0.378</b>	<b>0.439</b>	<b>0.658</b>	0.273	0.559	0.748	0.260
PolarFormer [22]	R101	900 $\times$ 1600	0.396	0.458	0.700	0.269	0.375	0.839	0.245
FCOS3D [49]	R101	900 $\times$ 1600	0.343	0.415	0.725	0.263	0.422	1.292	<b>0.153</b>
PGD [48]	R101	900 $\times$ 1600	0.369	0.428	0.683	<b>0.260</b>	0.439	1.268	0.185
DETR3D [52]	R101	900 $\times$ 1600	0.346	0.425	0.773	0.268	0.383	0.842	0.216
PETR [32]	R101	512 $\times$ 1408	0.357	0.421	0.710	0.270	0.490	0.885	0.224
BEVDet [20]	SwinT	512 $\times$ 1408	0.349	0.417	0.637	0.269	0.490	0.914	0.268
BEVDepth [27]	R101	512 $\times$ 1408	0.376	0.408	0.659	0.267	0.543	1.059	0.335
UniVision	R101	512 $\times$ 1408	<b>0.413</b>	<b>0.490</b>	<b>0.600</b>	0.263	<b>0.366</b>	<b>0.731</b>	0.211
PETRv2 [33]	R101	640 $\times$ 1600	0.421	0.524	0.681	0.267	0.357	<b>0.377</b>	<b>0.186</b>
PolarFormer [22]	R101	900 $\times$ 1600	0.432	0.528	0.648	0.270	0.348	0.409	0.201
BEVFormer [28]	R101	900 $\times$ 1600	0.416	0.517	0.673	0.274	0.372	0.394	0.198
BEVDet4D [19]	Swin-B	640 $\times$ 1600	0.396	0.515	0.619	<b>0.260</b>	0.361	0.399	0.189
UniVision4D (2frame)	R101	512 $\times$ 1408	0.439	0.535	0.580	0.264	<b>0.310</b>	0.491	0.200
UniVision4D (4frame)	R101	512 $\times$ 1408	<b>0.452</b>	<b>0.546</b>	<b>0.556</b>	0.264	0.330	0.451	0.199

Table 2. Results on nuScenes detection benchmark (val). Methods with \* are trained with aligned training settings, including input resolution, backbone, batch size, learning rate, etc. for a fair comparison.

when we use the same training settings for a fair comparison, UniVision shows superiority over other methods [20, 27, 28]. Specifically, UniVision achieves 2.4 points gain in mAP and 1.1 points gain in NDS against BEVDepth with the 512  $\times$  1408 image resolution. When we scale up the model and incorporate UniVision with temporal inputs, it further outperforms SOTA temporal-based detectors by a remarkable margin. UniVision achieves this with a smaller input resolution and it does not use the CBGS [65].

**OpenOccupancy.** The results on the OpenOccupancy benchmark are shown in Tab. 3. UniVision significantly surpasses recent vision-based occupancy methods including MonoScene [3], TPVFormer [21] and C-CONet [51] by 7.3 points, 6.5 points and 1.5 points in mIoU, respectively. Also, UniVision surpasses some LiDAR-based methods like LMSCNet [41] and JS3C-Net [55].

**Occ3D.** Tab. 4 lists the results on the Occ3D benchmark. With different input image resolutions, UniVision significantly

outperforms recent vision-based methods [20, 28, 53] by more than 2.7 points and 1.8 points in mIoU. Notably, BEVFormer and BEVDet-stereo load pre-trained weights and use temporal input in inference, while UniVision does not use them but still achieves better performance.

#### 4.4. Ablation Studies

**Effectiveness of components in detection task.** We show the ablation study for the detection task in Tab. 5. When we insert the BEV-based global feature extraction branch into the baseline model, the performance is boosted by 1.7 mAP and 3.0 NDS. When we add voxel-based occupancy task as an auxiliary task to the detector, the model has an improvement of 1.6 points gain in mAP. When we explicitly introduce the cross-representation interaction from the voxel features, the model achieves the best performance, which is improved by 3.5 points in mAP and 4.2 points in NDS compared with the baseline.

Method	Modality	mIoU	barrier	bycycle	bus	car	cons. vehi.	motorcycle	pedestrian	traffic_cone	trailer	truck	driv. surf.	other flat	sidewalk	terrain	manmade	vegetation
			■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■
LMSCNet [41]	LiDAR	11.5	12.4	4.2	12.8	12.1	6.2	4.7	6.2	6.3	8.8	7.2	24.2	12.3	16.6	14.1	13.9	22.2
JS3C-Net [55]	LIDAR	12.5	14.2	3.4	13.6	12.0	7.2	4.3	7.3	6.8	9.2	9.1	27.9	15.3	14.9	16.2	14.0	24.9
L-CONet [51]	LiDAR	15.8	17.5	5.2	13.3	18.1	7.8	5.4	9.6	5.6	13.2	13.6	34.9	21.5	22.4	21.7	19.2	23.5
AICNet [26]	Camera & Depth	10.6	11.5	4.0	11.8	12.3	5.1	3.8	6.2	6.0	8.2	7.5	24.1	13.0	12.8	11.5	11.6	20.2
3DSketch [7]	Camera & Depth	10.7	12.0	5.1	10.7	12.4	6.5	4.0	5.0	6.3	8.0	7.2	21.8	14.8	13.0	11.8	12.0	21.2
MonoScene [3]	Camera	6.9	7.1	3.9	9.3	7.2	5.6	3.0	5.9	4.4	4.9	4.2	14.9	6.3	7.9	7.4	<b>10.0</b>	7.6
TPVFormer [21]	Camera	7.8	9.3	4.1	11.3	10.1	5.2	4.3	5.9	5.3	<b>6.8</b>	6.5	13.6	9.0	8.3	8.0	9.2	8.2
C-CONet [51]	Camera	12.8	13.2	8.1	<b>15.4</b>	17.2	6.3	11.2	10.0	8.3	4.7	12.1	<b>31.4</b>	18.8	18.7	16.3	4.8	8.2
UniVision	Camera	<b>14.3</b>	<b>14.7</b>	<b>9.6</b>	12.9	<b>17.2</b>	<b>8.6</b>	<b>12.1</b>	<b>12.1</b>	<b>9.1</b>	6.3	<b>13.0</b>	30.5	<b>23.0</b>	<b>19.9</b>	<b>18.0</b>	9.1	<b>13.3</b>

Table 3. Results on OpenOccupancy benchmark. UniVision achieves state-of-the-art performance among camera-based methods and has comparable or better results than some of the LiDAR-based methods.

Method	Resolution	mIoU	others	barrier	bycycle	bus	car	cons. vehi.	motorcycle	pedestrian	traffic_cone	trailer	truck	driv. surf.	other flat	sidewalk	terrain	manmade	vegetation
			■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■
BEVFormer [28]	256 × 704	29.2	4.8	35.0	0.0	35.5	41.4	10.4	1.8	17.5	10.1	23.3	26.7	79.3	36.6	47.7	51.9	39.6	35.0
SurroundOcc [53]	256 × 704	31.3	9.2	33.7	17.5	34.6	39.9	16.7	18.4	22.4	20.8	25.6	27.3	75.4	35.1	43.2	47.1	34.5	30.2
OpenOccupancy [51]	256 × 704	32.5	9.7	40.2	18.8	36.6	44.1	8.1	19.5	22.9	24.2	24.5	28.9	77.6	37.3	45.2	49.2	35.5	30.6
BEVDet-depth [20]	256 × 704	31.4	5.9	38.7	0.5	38.6	44.9	14.4	8.4	17.0	14.6	24.1	29.5	79.7	39.5	49.6	53.1	39.8	34.8
BEVDet-stereo [20]	256 × 704	34.8	8.1	42.0	5.7	41.3	47.7	21.0	15.5	19.4	17.7	29.8	34.7	<b>80.4</b>	38.7	<b>51.8</b>	<b>55.1</b>	<b>44.4</b>	<b>38.8</b>
UniVision	256 × 704	<b>37.5</b>	<b>11.0</b>	<b>44.7</b>	<b>23.1</b>	<b>43.0</b>	<b>50.5</b>	<b>21.6</b>	<b>24.9</b>	<b>26.9</b>	<b>25.7</b>	<b>30.7</b>	<b>35.8</b>	79.8	<b>41.4</b>	49.1	53.8	40.3	34.7
BEVFormer [28]	512 × 1408	34.7	7.1	40.7	9.8	40.1	47.8	16.3	17.5	23.3	20.7	27.4	33.2	81.3	39.7	50.4	54.3	43.2	36.8
SurroundOcc [53]	512 × 1408	34.4	8.7	39.2	19.7	41.4	46.2	18.7	20.6	26.4	23.3	27.0	32.5	78.0	38.3	46.6	49.6	36.7	31.6
OpenOccupancy [51]	512 × 1408	36.1	10.4	45.7	23.6	42.4	49.3	14.8	24.6	27.7	27.8	27.6	33.3	79.2	39.8	47.1	50.5	37.7	31.8
BEVDet-depth [20]	512 × 1408	33.7	6.6	41.2	7.0	42.7	48.4	18.4	12.9	22.0	18.2	28.2	33.2	80.1	39.7	49.1	52.1	39.9	33.8
BEVDet-stereo [20]	512 × 1408	38.0	8.6	45.9	14.3	<b>46.0</b>	51.2	<b>23.8</b>	18.9	24.1	22.3	<b>33.6</b>	37.9	<b>81.5</b>	40.5	<b>52.6</b>	<b>55.9</b>	<b>46.9</b>	<b>41.2</b>
UniVision	512 × 1408	<b>39.8</b>	<b>11.3</b>	<b>47.1</b>	<b>27.6</b>	45.8	<b>54.2</b>	22.9	<b>28.6</b>	<b>31.0</b>	<b>28.7</b>	31.8	<b>38.5</b>	81.2	<b>42.5</b>	51.3	54.7	41.7	36.6

Table 4. Results on Occ3D benchmark. UniVision achieves state-of-the-art performance with different input image resolutions. The results of the compared methods are reproduced with their official code base. \* Note that BEVFormer uses video input and BEVDet-stereo uses depth pre-training and stereo input but UniVision does not.

**Effectiveness of components in occupancy task.** We show the ablation study for the occupancy task in Tab. 6. The voxel-based local feature extraction network brings an improvement of 1.96 mIoU gain to the baseline model. When the detection task is introduced as an auxiliary supervision signal, the model performance is boosted by 0.4 in mIoU. When we explicitly fuse the glocal BEV features with the voxel features, the model achieves the best performance with 2.64 points gain in mIoU compared with the baseline.

**How do detection task and occupancy task influence each other?** Tab. 5 and Tab. 6 show that both the detection task and the occupancy task benefit each other in our UniVision framework. For the detection task, the occupancy supervision can improve the mAP and mATE metrics, which indicates that voxel-wise semantic learning effectively improves the detector’s awareness of object geometry, *i.e.*, centerness and scale. For the occupancy task, the detection supervision significantly improves the performance for foreground categories, *i.e.*, detection categories, thus resulting in an overall improvement.

**Effectiveness of joint Occ-Det augmentation, Ex-Im view transform and progressive loss weight adjustment.** We show the effectiveness of the joint Occ-Det spatial augmentation, the Ex-Im view transform module, and the progressive loss weight adjustment strategy in Tab. 7. It shows significant improvements in the detection task and the occupancy task with the proposed spatial augmentation and the proposed view transform module on the mIoU, mAP, and NDS metrics. The loss weight adjustment strategy enables the efficient training of the multi-task framework. Without this, the training of the unified framework cannot converge and the performance is very low.

#### 4.5. Qualitative Results

We show qualitative results of UniVision in Fig. 4, which include the detection results on the 2D image plane, the detection results on the BEV plane, and the corresponding occupancy prediction results. The UniVision framework can simultaneously produce high-quality prediction results for both 3D detection and occupancy prediction tasks with a unified network.

Global Branch	Occ Sup.	Voxel Inter.	mAP $\uparrow$	NDS $\uparrow$	mATE $\downarrow$	mASE $\downarrow$	mAOE $\downarrow$	mAVE $\downarrow$	mAAE $\downarrow$
×	×	×	0.271	0.352	0.725	0.287	0.663	0.911	0.256
✓	×	×	0.288	0.382	0.718	0.291	0.623	<b>0.776</b>	0.214
✓	✓	×	0.304	0.384	0.707	0.285	0.680	0.789	0.222
✓	×	✓	0.297	0.389	0.684	0.285	0.615	0.804	<b>0.197</b>
✓	✓	✓	<b>0.306</b>	<b>0.394</b>	<b>0.662</b>	<b>0.284</b>	<b>0.601</b>	0.840	0.205

Table 5. Ablation study for the detection task. *Occ Sup.* denotes adding the voxel branch and using the occupancy task as auxiliary supervision without interaction. *Voxel Inter.* denotes explicitly using the voxel features for cross-representation interaction.



Figure 4. Qualitative results of UniVision framework, including the detection results on the 2D image plane, the detection results on the BEV plane, and the corresponding occupancy prediction results.

Local Branch	Det Sup.	BEV Inter.	mIoU	mIoU <sub>fore</sub>	mIoU <sub>back</sub>
×	×	×	34.58	29.42	47.13
✓	×	×	36.52	31.15	49.66
✓	✓	×	36.96	32.22	49.18
✓	×	✓	36.51	30.93	<b>49.99</b>
✓	✓	✓	<b>37.12</b>	<b>32.37</b>	49.38

Table 6. Ablation study for the occupancy task. *Det Sup.* denotes adding the BEV branch and using the detection task as auxiliary supervision without interaction. *BEV Inter.* denotes explicitly using the BEV features for cross-representation interaction.  $mIoU_{fore}$  denotes the mIoU of the foreground object classes.  $mIoU_{back}$  denotes the mIoU of the background object classes.

Occ-Det Aug.	Ex-Im Trans.	Prog. Adj.	mIoU	mAP	NDS
×	-	-	36.44	0.271	0.372
✓	-	-	<b>37.12</b>	<b>0.306</b>	<b>0.394</b>
-	×	-	36.96	0.304	0.384
-	✓	-	<b>37.46</b>	<b>0.312</b>	<b>0.397</b>
-	-	×	-	-	-
-	-	✓	<b>37.12</b>	<b>0.306</b>	<b>0.394</b>

Table 7. Ablation study for joint Occ-Det augmentation, Ex-Im view transform and progressive loss weight adjustment. – denotes that the model can not converge and the performance is very low.

## 5. Conclusion

In this paper, we present UniVision, a unified framework for vision-centric 3D perception that can simultaneously handle the occupancy prediction task and the 3D detection task. To achieve this, we propose the explicit-implicit view transform module for complementary 2D-3D feature transformation. We propose a local-global feature extraction and fusion module for efficient and adaptive multi-representation feature extraction, enhancement, and interaction. Besides, we propose a joint occ-det data augmentation strategy and a progressive loss weight adjustment strategy

for efficient and stable multi-task framework training. UniVision achieves state-of-the-art performance on four benchmarks, including nuScenes LiDAR segmentation, nuScenes detection, OpenOccupancy, and Occ3D.

## References

- [1] Garrick Brazil and Xiaoming Liu. M3d-rpn: Monocular 3d region proposal network for object detection. In *Int. Conf. Comput. Vis.*, 2019. 2
- [2] Holger Caesar, Varun Bankiti, Alex H Lang, and et al. nuscenes: A multimodal dataset for autonomous driving. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020. 2, 6
- [3] Anh-Quan Cao and Raoul de Charette. Monoscene: Monocular 3d semantic scene completion. *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021. 1, 2, 7, 8
- [4] Thomas Cashman, Corey Vogel, and Paul Newman. Contextual scene completion in 3d rgb-d indoor scenes. In *International Conference on Robotics and Automation (ICRA)*, 2018. 2
- [5] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2018. 2
- [6] Xiaozhi Chen, Kaustav Kundu, Ziyu Zhang, Huimin Ma, Sanja Fidler, and Raquel Urtasun. Monocular 3d object detection for autonomous driving. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2147–2156, 2016. 2
- [7] Xiaokang Chen, Kwan-Yee Lin, Chen Qian, Gang Zeng, and Hongsheng Li. 3d sketch-aware semantic scene completion via semi-supervised structure prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4193–4202, 2020. 8
- [8] Zehui Chen, Zhenyu Li, Shiquan Zhang, Liangji Fang, Qin-hong Jiang, and Feng Zhao. Bevdistill: Cross-modal bev distillation for multi-view 3d object detection. 2022. 1
- [9] Zhiyu Chong, Xinzhu Ma, Hong Zhang, Yuxin Yue, Haojie Li, Zhihui Wang, and Wanli Ouyang. Monodistill: Learning spatial features for monocular 3d object detection, 2022. 1
- [10] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. 2017. 1
- [11] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. Centernet: Keypoint triplets for object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6569–6578, 2019. 2
- [12] 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. 2021. 2
- [13] Saurabh Gupta, Pablo Arbeláez, Ross Girshick, and Jitendra Malik. Learning rich features from rgb-d images for object detection and segmentation. In *European Conference on Computer Vision (ECCV)*, 2014. 2
- [14] Chunrui Han, Jianjian Sun, Zheng Ge, Jinrong Yang, Runpei Dong, Hongyu Zhou, Weixin Mao, Yuang Peng, and Xiangyu Zhang. Exploring recurrent long-term temporal fusion for multi-view 3d perception. 2023. 1, 2
- [15] Ali Hassani, Steven Walton, Jiachen Li, Shen Li, and Humphrey Shi. Neighborhood attention transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6185–6194, June 2023. 4
- [16] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross B. Girshick. Mask R-CNN. In *International Conference on Computer Vision (ICCV)*, 2017. 2
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, et al. Deep residual learning for image recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2016. 4, 6
- [18] Yu Hong, Hang Dai, and Yong Ding. Cross-modality knowledge distillation network for monocular 3d object detection. pages 87–104, Nov. 2022. 1
- [19] Junjie Huang and Guan Huang. Bevdet4d: Exploit temporal cues in multi-camera 3d object detection. *arXiv preprint arXiv:2203.17054*, 2022. 7
- [20] 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. 1, 2, 3, 5, 6, 7, 8
- [21] Yuanhui Huang, Wenzhao Zheng, Yunpeng Zhang, Jie Zhou, and Jiwen Lu. Tri-perspective view for vision-based 3d semantic occupancy prediction. *arXiv preprint arXiv:2302.07817*, 2023. 1, 2, 6, 7, 8
- [22] 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. In *AAAI*, 2023. 7
- [23] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [24] Iasonas Kokkinos. Ubernet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory. *arXiv preprint arXiv:1611.06612*, 2017. 2
- [25] Alex Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *CVPR*, 2019. 2
- [26] Jie Li, Kai Han, Peng Wang, Yu Liu, and Xia Yuan. Anisotropic convolutional networks for 3d semantic scene completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3351–3359, 2020. 8
- [27] 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. 4, 6, 7
- [28] Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Qiao Yu, and Jifeng Dai. Bevformer: Learning bird’s-eye-view representation from multi-camera images via spatiotemporal transformers. *arXiv preprint arXiv:2203.17270*, 2022. 1, 2, 6, 7, 8
- [29] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. 4

- [30] Xuewu Lin, Tianwei Lin, Zixiang Pei, Lichao Huang, and Zhizhong Su. Sparse4d: Multi-view 3d object detection with sparse spatial-temporal fusion. *CoRR*, abs/2211.10581, 2022. 2
- [31] Haisong Liu, Yao Teng, Tao Lu, Haiguang Wang, and Limin Wang. Sparsebev: High-performance sparse 3d object detection from multi-camera videos. *CoRR*, abs/2308.09244, 2023. 2
- [32] Yingfei Liu, Tiancai Wang, Xiangyu Zhang, and Jian Sun. Petr: Position embedding transformation for multi-view 3d object detection. *arXiv preprint arXiv:2203.05625*, 2022. 2, 7
- [33] Yingfei Liu, Junjie Yan, Fan Jia, Shuailin Li, Qi Gao, Tiancai Wang, Xiangyu Zhang, and Jian Sun. Petrv2: A unified framework for 3d perception from multi-camera images. *arXiv preprint arXiv:2206.01256*, 2022. 7
- [34] 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 *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. 1
- [35] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 6
- [36] Xinzhu Ma, Yinmin Zhang, Dan Xu, et al. Delving into localization errors for monocular 3d object detection. In *CVPR*, 2021. 2
- [37] Ishan Misra, Abhinav Shrivastava, Abhinav Gupta, and Judy He. Cross-stitch networks for multi-task learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2
- [38] Dennis Park, Rares Ambrus, and Vitor others Guizilini. Is pseudo-lidar needed for monocular 3d object detection? In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021. 2
- [39] Jonah Philion and Sanja Fidler. Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d. In *European Conference on Computer Vision*, pages 194–210. Springer, 2020. 2, 3
- [40] Cody Reading, Ali Harakeh, Julia Chae, and Steven L. Waslander. Categorical depth distributionnetwork for monocular 3d object detection. *CVPR*, 2021. 1, 2
- [41] Luis Roldao, Raoul de Charette, and Anne Verroust-Blondet. Lmscnet: Lightweight multiscale 3d semantic completion. In *2020 International Conference on 3D Vision (3DV)*, pages 111–119. IEEE, 2020. 7, 8
- [42] Shuran Song, Fisher Yu, Andy Zeng, Angel X Chang, Manolis Savva, and Thomas Funkhouser. Semantic scene completion from a single depth image. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [43] Xiaoyu Tian, Tao Jiang, Longfei Yun, Yue Wang, Yilun Wang, and Hang Zhao. Occ3d: A large-scale 3d occupancy prediction benchmark for autonomous driving. *arXiv preprint arXiv:2304.14365*, 2023. 1, 2, 6
- [44] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019. 2
- [45] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3
- [46] Chengliang Wang, Haojian Ning, Xinrun Chen, and Shiyong Li. Db-unet: Mlp based dual branch unet for accurate vessel segmentation in octa images. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023. 6, 7
- [47] Shihao Wang, Yingfei Liu, Tiancai Wang, Ying Li, and Xiangyu Zhang. Exploring object-centric temporal modeling for efficient multi-view 3d object detection. 2023. 1, 2
- [48] 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. 7
- [49] 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*, pages 913–922, 2021. 2, 7
- [50] Wenhai Wang, Jifeng Dai, Zhe Chen, Zhenhang Huang, Zhiqi Li, Xizhou Zhu, Xiaowei Hu, Tong Lu, Lewei Lu, Hongsheng Li, et al. Internimage: Exploring large-scale vision foundation models with deformable convolutions. *arXiv preprint arXiv:2211.05778*, 2022. 4
- [51] Xiaofeng Wang, Zheng Zhu, Wenbo Xu, Yunpeng Zhang, Yi Wei, Xu Chi, Yun Ye, Dalong Du, Jiwen Lu, and Xingang Wang. Openoccupancy: A large scale benchmark for surrounding semantic occupancy perception. *arXiv preprint arXiv:2303.03991*, 2023. 1, 2, 4, 6, 7, 8
- [52] 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. 2, 7
- [53] Yi Wei, Linqing Zhao, Wenzhao Zheng, Zheng Zhu, Jie Zhou, and Jiwen Lu. Surroundocc: Multi-camera 3d occupancy prediction for autonomous driving. *arXiv preprint arXiv:2303.09551*, 2023. 1, 2, 6, 7, 8
- [54] Jianxiong Xiao, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Scene completion using millions of photographs. *ACM Transactions on Graphics (TOG)*, 2010. 2
- [55] Xu Yan, Jiantao Gao, Jie Li, Ruimao Zhang, Zhen Li, Rui Huang, and Shuguang Cui. Sparse single sweep lidar point cloud segmentation via learning contextual shape priors from scene completion. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 3101–3109, 2021. 7, 8
- [56] Yan Yan, Yuxing Mao, and Bo Li. SECOND: sparsely embedded convolutional detection. *Sensors*, 2018. 2, 4
- [57] Dongqiangzi Ye, Zixiang Zhou, Weijia Chen, Yufei Xie, Yu Wang, Panqu Wang, and Hassan Foroosh. Lidarmultinet: Towards a unified multi-task network for lidar perception. *arXiv preprint arXiv:2209.09385*, 2022. 2, 7
- [58] Tianwei Yin, Xingyi Zhou, and Philipp Krähenbühl. Center-based 3d object detection and tracking. *CVPR*, 2021. 2, 4

- [59] Renrui Zhang, Han Qiu, Tai Wang, Xuanzhuo Xu, Ziyu Guo, Yu Qiao, Peng Gao, and Hongsheng Li. Monodetr: Depth-aware transformer for monocular 3d object detection. *arXiv preprint arXiv:2203.13310*, 2022. 2
- [60] Yunpeng Zhang, Jiwen Lu, and Jie Zhou. Objects are different: Flexible monocular 3d object detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021. 2
- [61] Yang Zhang, Zixiang Zhou, Philip David, Xiangyu Yue, Zerong Xi, Boqing Gong, and Hassan Foroosh. Polarnet: An improved grid representation for online lidar point clouds semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9601–9610, 2020. 6, 7
- [62] Yunpeng Zhang, Zheng Zhu, and Dalong Du. Occformer: Dual-path transformer for vision-based 3d semantic occupancy prediction. *arXiv preprint arXiv:2304.05316*, 2023. 1, 2, 6, 7
- [63] Zhi Zhang, Chen Cao, Chengrui Zhang, Shifeng Li, and Nanning Zheng. Multi-task learning with low-level tasks for instance segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018. 2
- [64] Hui Zhou, Xinge Zhu, Xiao Song, Yuexin Ma, Zhe Wang, Hongsheng Li, and Dahua Lin. Cylinder3d: An effective 3d framework for driving-scene lidar semantic segmentation. *arXiv preprint arXiv:2008.01550*, 2020. 7
- [65] 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. 6, 7
- [66] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 3
- [67] Zhuofan Zong, Dongzhi Jiang, Guanglu Song, Zeyue Xue, Jingyong Su, Hongsheng Li, and Yu Liu. Temporal enhanced training of multi-view 3d object detector via historical object prediction. 2023. 1, 2