# SurroundOcc: Multi-Camera 3D Occupancy Prediction for Autonomous Driving

Yi Wei[1,2]*, Linqing Zhao[3]*, Wenzhao Zheng[1,2], Zheng Zhu[4], Jie Zhou[1,2], Jiwen Lu[1,2]†

[1]Beijing National Research Center for Information Science and Technology, China
[2]Department of Automation, Tsinghua University, China
[3]School of Electrical and Information Engineering, Tianjin University, China
[4]PhiGent Robotics

{y-wei19,zhengwz18}@mails.tsinghua.edu.cn; linqingzhao@tju.edu.cn;
zhengzhu@ieee.org; {jzhou, lujiwen}@tsinghua.edu.cn

## Abstract

*3D scene understanding plays a vital role in vision-based autonomous driving. While most existing methods focus on 3D object detection, they have difficulty describing real-world objects of arbitrary shapes and infinite classes. Towards a more comprehensive perception of a 3D scene, in this paper, we propose a SurroundOcc method to predict the 3D occupancy with multi-camera images. We first extract multi-scale features for each image and adopt spatial 2D-3D attention to lift them to the 3D volume space. Then we apply 3D convolutions to progressively upsample the volume features and impose supervision on multiple levels. To obtain dense occupancy prediction, we design a pipeline to generate dense occupancy ground truth without expansive occupancy annotations. Specifically, we fuse multi-frame LiDAR scans of dynamic objects and static scenes separately. Then we adopt Poisson Reconstruction to fill the holes and voxelize the mesh to get dense occupancy labels. Extensive experiments on nuScenes and SemanticKITTI datasets demonstrate the superiority of our method. Code and dataset are available at https://github.com/weiyithu/SurroundOcc.*

## 1. Introduction

Understanding the 3D geometry of the surrounding scene serves as the basic procedure in an autonomous driving system. While LiDAR is a direct solution to capture this geometric information, it suffers from high-cost sensors and sparse scanned points, limiting its further application. Recently, vision-centric autonomous driving has attracted extensive attention as a promising direction. Taking multi-
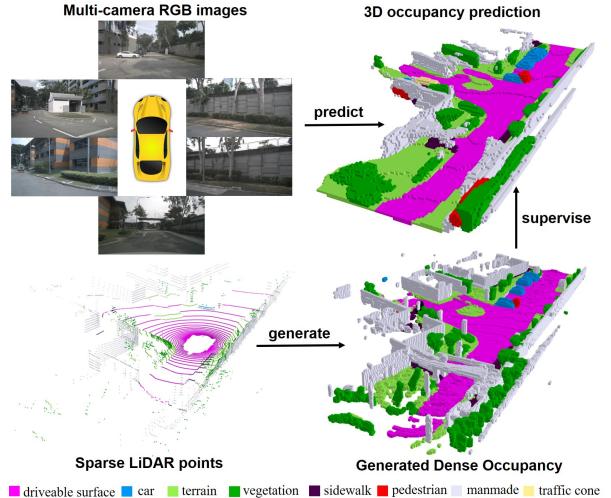


Figure 1. The overview of SurroundOcc. Given multi-camera images, our method can predict volumetric occupancy of surrounding 3D scenes. To train the network, we design a pipeline to generate dense occupancy labels with sparse LiDAR points. **Better viewed when zoomed in.**

camera images as inputs, it has demonstrated competitive performance in various 3D perception tasks including depth estimation [18, 57], 3D object detection [29, 28, 35, 33, 21], and semantic map construction [65, 20, 1].

Although multi-camera 3D object detection plays an important role [29, 28, 33, 21] in vision-based 3D perception, it is easy to suffer from the long-tail problem and difficult to recognize all classes of objects in the real world. Complementary to 3D detection, reconstructing surrounding 3D scenes can better help the downstream perception tasks. Recent works [18, 57] incorporate information from multiple views and predict surrounding depth maps. However, depth maps only predict the nearest occupied point in each optical ray and are unable to recover the occluded parts of the 3D

---

*Equal contribution.
†Corresponding author.

scene. Different from depth-based methods, another trend [8, 22] is to directly predict the 3D occupancy of the scene. It describes a 3D scene by assigning an occupied probability to each voxel in the 3D space. We advocate 3D occupancy to be a good 3D representation for multi-camera scene reconstruction, which naturally guarantees the multi-camera geometry consistency and is able to recover occluded parts. Also, it is flexible to extend to other 3D downstream tasks such as 3D semantic segmentation [68, 72, 16]. As one of the pioneering works, MonoScene [8] infers the dense 3D voxelized semantic scene with monocular images. However, simply fusing multi-camera results with cross-camera post-processing will lead to low performance [29]. TPVFormer [22] uses sparse LiDAR points as supervision, which results in sparse occupancy prediction.

To address this, we propose a SurroundOcc method, which aims to predict dense and accurate 3D occupancy with multi-camera images input. We first use a 2D backbone network to extract multi-scale feature maps from each image. Then we perform 2D-3D spatial attention to lift multi-camera image information to 3D volume features instead of BEV features. A 3D convolution network is then employed to progressively upsample the low-resolution volume features and fuse them with high-resolution ones to obtain fine-grained 3D representations. At each level, we use a decayed weighted loss to supervise the network. To get dense predictions, we need dense occupancy labels. However, the mainstream multi-camera dataset nuScenes [7] only provides sparse LiDAR points. To avoid expensive occupancy annotations, we devise a pipeline to generate dense occupancy ground truth only with the existing 3D detection and 3D semantic segmentation labels. Specifically, we first combine multi-frame points of dynamic objects and static scenes respectively. Then we leverage Poisson Reconstruction [24] algorithm to further fill the holes. Finally, NN and voxelization are used to obtain dense 3D occupancy labels.

With the dense occupancy ground truth, we train the model and compare it with other state-of-the-art methods on nuScenes [7] dataset. Both the quantitative results and visualizations demonstrate the effectiveness of our method. Moreover, we further conduct experiments on SemanticKITTI dataset [2]. Although our method is not designed for the monocular setting, it achieves state-of-the-art performance on the monocular 3D semantic scene completion benchmark.

## 2. Related Work

**Voxel-based Scene Representation:** How to effectively represent a 3D scene lies at the core of autonomous driving perception. Voxel-based scene representation voxelizes the 3D space into discretized voxels and describes each voxel by a vector feature. It has empowered the success of numerous methods on the lidar segmentation [32, 51, 12, 61, 60]

and 3D scene completion [8, 45, 10, 26, 58, **?**, **?**] tasks. For the 3D occupancy prediction task, we also advocate the voxel representation as it is more suitable to model the occupancy field of a 3D scene. MonoScene [8] is the first work to reconstruct outdoor scenes using only RGB inputs. TPVFormer [22] further generalizes it to multi-camera 3D semantic occupancy prediction. However, its lack of dense supervision results in sparse occupancy prediction. Differently, we devise a pipeline to generate dense occupancy ground truth for training and our occupancy prediction is much denser.

**3D Scene Reconstruction:** 3D reconstruction [41, 48, 47, 40, 70, 17, 5, 37, 49, 6] is a traditional but important topic in computer vision. One way is through depth estimation which predicts a depth value for each pixel in the image. While early methods require full depth annotations to supervise the depth estimation model [15, 25, 66], later research focuses on self-supervised depth estimation as it does not require intensive human annotations [70, 17, 5, 67, 62, 56, 69, 43, 9]. Recently, SurroundDepth [57] further incorporates the interactions between surrounding views to capture more spatial correlations. Different from depth estimation, 3D scene reconstruction methods [23, 37, 49, 6, 8] directly reconstruct a comprehensive and accurate 3D geometry of a scene. SurfaceNet [23] employs a 3D convolutional network to transform RGB colors to 3D surface occupancy from two images. Atlas [37] further extends it to the multi-view setting and utilized learned features to predict occupancy. NeuralRecon [49] and TransformerFusion [6] fuse the learned image features from different views in an online manner for more accurate 3D reconstructions. However, most of these 3D scene reconstruction methods are designed for indoor scenes, which are different from the multi-camera setting in the outdoor environment.

**Vision-based 3D Perception:** The lack of direct geometric input demands vision-based 3D surround perception methods to infer the 3D scene geometry based on semantic cues. Depth-based methods explicitly predict depth maps for image inputs to extract 3D geometric information of the scene before perception [54, 36, 44, 39, 28, 35, 30, 42, 21, 65]. The widely adopted pipeline is to predict categorical depth distributions and leverage them to project semantic features into 3D space [42]. Other methods implicitly learn 3D features without producing explicit depth maps [53, 29, 55, 34, 52, 63]. For example, BEVFormer [29] adopts cross attention to progressively refine BEV grid features from 2D image features. While most existing works employ BEV representation to describe a scene, we propose to reconstruct a 3D scene using volumetric occupancy representation, which provides the more fine-grained and comprehensive modeling of the scene.
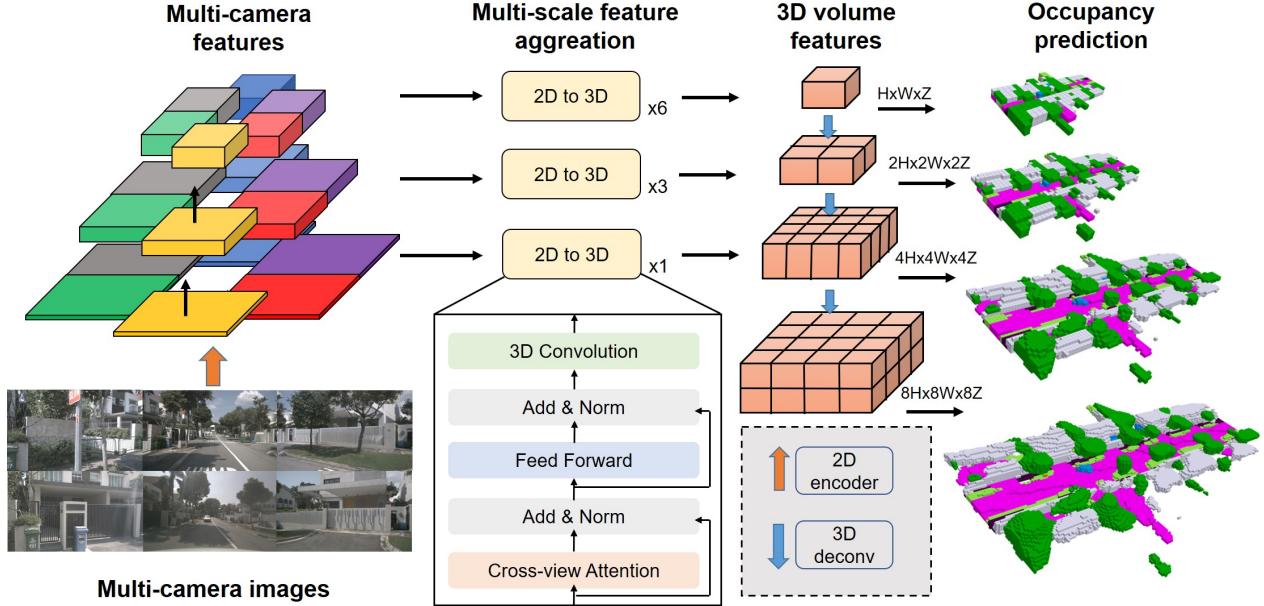
Figure 2. The pipeline of the proposed method. First, we use a backbone to extract multi-scale features of multi-camera images. Then we adopt 2D-3D spatial attention to fuse multi-camera information and construct 3D volume features in a multi-scale fashion. Finally, the 3D deconvolution layer is used to upsample 3D volumes and occupancy prediction is supervised in each level.
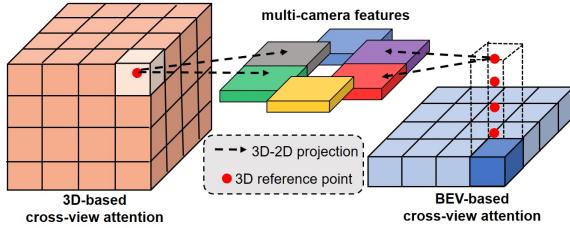


Figure 3. The comparison of 3D-based and BEV-based cross-view attention. The 3D-based attention can better preseve 3D information. For each 3D volume query, we project it to the corresponding 2D views to sample features.

## 3. Approach

### 3.1. Problem Formulation

In this work, we aim to predict the 3D occupancy of surrounding scenes with multi-camera images $I = \{I^1, I^2, \cdots I^N\}$. Formally, the 3D occupancy predcition is represented as:

$$V = G(I^1, I^2, \cdots I^N) \qquad (1)$$

where $G$ is an neural network and $V \in \mathbb{R}^{H \times W \times Z}$ is the 3D occupancy. The value of $V$ is between 0 and 1, representing the occupied probability of the grids. Lifting $V$ to an $(L, H, W, Z)$ tensor, we can obtain the 3D semantic occupancy, where $L$ is the class number and class 0 means non-occupied grids.

3D occupancy is a good representation for multi-camera 3D scene reconstruction. First, since 3D occupancy is

predicted in 3D space, it theoretically satisfies the multi-camera consistency. Second, it is possible for networks to predict occluded areas according to the surrounding semantic information, which is unavailable in depth estimation. Third, 3D occupancy is easy to extend to other downstream tasks, such as 3D semantic segmentation and scene flow estimation.

### 3.2. Overview

Figure 2 shows the pipeline of our method. Given a set of surrounding multi-camera images, we first use a backbone network (*e.g.* ResNet-101 [19]) to extract $N$ cameras' and $M$ levels' multi-scale features $X = \{\{X_i^j\}_{i=1}^N\}_{j=1}^M$. For each level, we use a transformer to fuse multi-camera features with spatial cross attention. The output of the 2D-3D spatial attention layer is a 3D volume feature instead of the BEV feature. Then the 3D convolution network is utilized to upsample and combine multi-scale volume features. The occupancy prediction in each level is supervised by the generated dense occupancy ground truth with a decayed loss weight.

### 3.3. 2D-3D Spatial Attention

Many 3D scene reconstruction methods [8, 37] integrate multi-view 2D features into 3D space by reprojecting 2D features back to the 3D volumes with known poses. The grid feature is calculated by simply averaging all 2D features in this grid. However, this kind of method assumes that different views contribute equally to the 3D volumes,
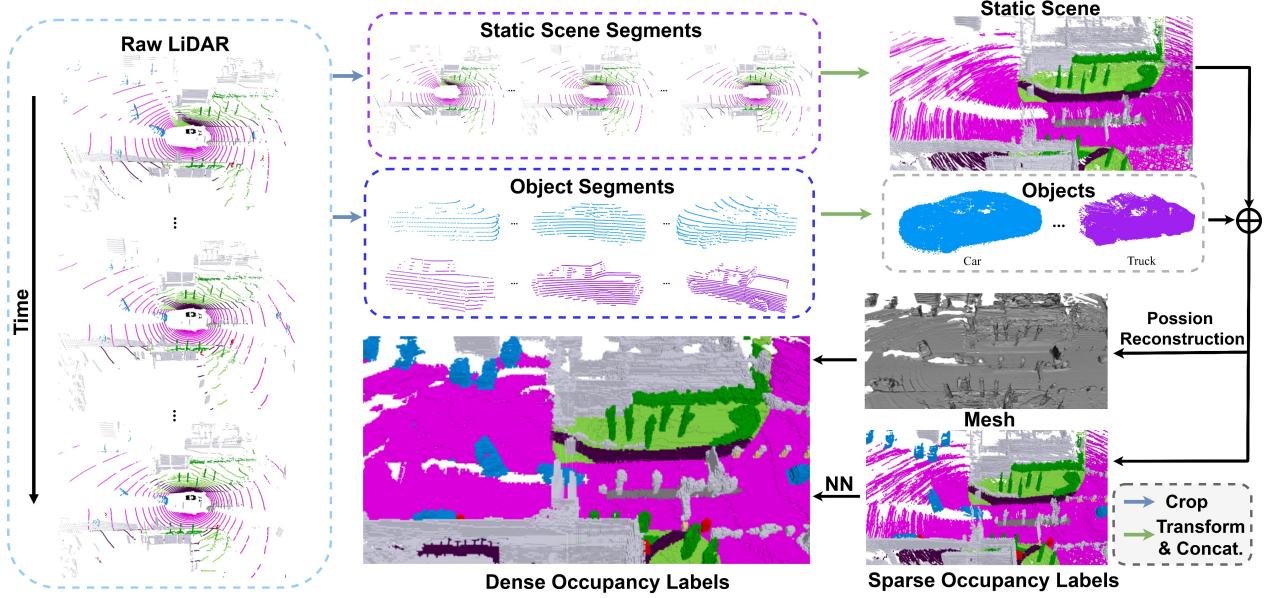
Figure 4. Dense occupancy ground truth generation. We first traverse all frames to stitch the multi-frame LiDAR points of dynamic objects and static scenes separately, and then merge them into a complete scene. Subsequently, we employ Poisson Reconstruction to densify the points and voxelize the resulting mesh to obtain a dense 3D occupancy. Finally, we use the Nearest Neighbor (NN) algorithm to assign semantic labels to dense voxels.

which does not always hold true, especially when some views are occluded or blurred.

To tackle the issue, we leverage cross-view attention to fuse multi-camera features. We project 3D reference points to 2D views and use deformable attention [71, 29] to query points and aggregate information. As shown in Figure 3, instead of 2D BEV queries, we build 3D volume queries to further reserve 3D space information. Specifically, 3D volume queries are defined as $Q \in \mathbb{R}^{C \times H \times W \times Z}$. For each query, we project its corresponding 3D point to 2D views according to the given intrinsic and extrinsic. We only use the views that the 3D reference point hits. Then we sample 2D features around these projected 2D positions. The output $F \in \mathbb{R}^{C \times H \times W \times Z}$ of the cross-view attention layer is a weighted sum of sampled features according to the deformable attention mechanism:

$$\text{DeformAttn}(q, p, x) = \sum_{i=1}^{N_{\text{head}}} \mathcal{W}_i \sum_{j=1}^{N_{\text{key}}} \mathcal{A}_{ij} \cdot \mathcal{W}'_i x(p + \Delta p_{ij})$$

$$F^p = \frac{1}{|\mathcal{V}_{\text{hit}}|} \sum_{i \in \mathcal{V}_{\text{hit}}} \text{DeformAttn}(Q^p, \mathcal{P}(q^p, i), X_i) \tag{2}$$

where $F^p$ and $Q^p$ indicate the $p$th element of the output features and 3D volume queries. $q^p$ is the corresponding 3D positions of queries and $P$ is the 3D to 2D project function. $\mathcal{V}_{\text{hit}}$ represents the hit views of 3D query points. $W_i$ and $W'_i$ are the learnable weights and $A_{ij} \in [0, 1]$ is the attention weight calculated by the dot product of query and key. $x(p + \Delta p_{ij})$ is the 2D feature at location $p + \Delta p_{ij}$. In-

stead of performing expensive 3D self-attention, we use the 3D convolution to interact features between neighboring 3D voxels.

### 3.4. Multi-scale Occupancy Prediction

We further extend the 2D-3D spatial attention to a multi-scale fashion. Different from 3D detection task, 3D scene reconstruction needs more low-level features to help the network learn fine-grained details. To tackle the issue, we design a 2D-3D U-Net architecture. Specifically, given multi-scale 2D features $\{\{X_i^j\}_{i=1}^N\}_{j=1}^M$, we adopt different number of 2D-3D spatial attention layers to extract multi-scale 3D volume features $\{F_j \in \mathbb{R}^{C_j \times H_j \times W_j \times Z_j}\}_{j=1}^M$. Then we upsample $j-1$th level 3D volume features $Y_{j-1}$ with 3D deconvolution layer and fuse it with $F_j$:

$$Y_j = F_j + \text{Deconv}(Y_{j-1}) \tag{3}$$

For each level, the network outputs an occupancy prediction result with different resolution $V_j \in \mathbb{R}^{C_j \times H_j \times W_j \times Z_j}$. To get powerful both high-level and low-level 3D features, the network is supervised at each scale. Specifically, we use the cross-entropy loss and scene-class affinity loss introduced in [8] as the supervision signals. For 3D semantic occupancy prediction, we adopt a multi-class cross-entropy loss and for 3D scene reconstruction we change it to a two-class formulation. Since the high-resolution prediction is more important, we use a decayed loss weight $\alpha_j = \frac{1}{2^j}$ for $j$th level supervision.

(a) RGB image      (b) single-frame LiDAR points

(c) Sparse occupancy labels      (d) Dense occupancy labels
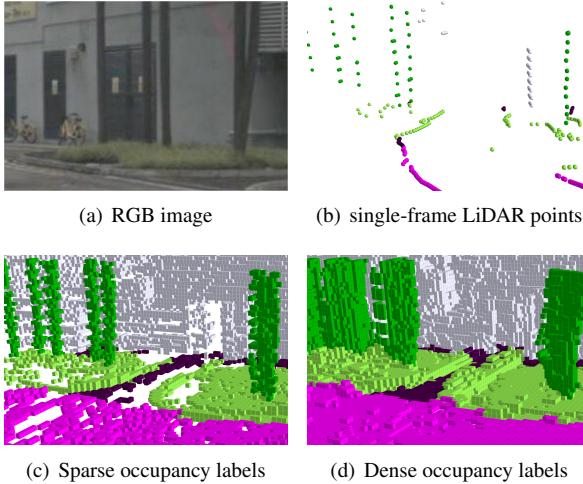
Figure 5. Comparison on different occupancy labels. Compared with single-frame LiDAR points and the sparse occupancy converted from multi-frame points, our dense voxels are able to provide more realistic occupancy labels.

## 4. Dense Occupancy Ground Truth

In our experiments, we find that the network supervised by sparse LiDAR points is unable to predict dense enough occupancy. Thus, it is necessary to generate dense occupancy labels. However, as mentioned in SemanticKITTI [2], it is complex and needs huge human efforts to annotate the dense occupancy of a 3D scene which has millions of voxels. To this end, we design a pipeline to generate dense occupancy ground truth leveraging existing 3D detection and 3D semantic segmentation labels without extra human annotations, which is shown in Figure 4. An intuitive way is to directly transform a multi-frame LiDAR point sequence into a unified coordinate system, and then voxelize the concatenated dense points into voxel grids. However, such a straightforward solution is only applicable to completely static scenes and ignores moving objects. Moreover, multi-frame point clouds are not dense enough and there still exist many holes, which will result in wrong occupancy labels. To address these issues, we propose to stitch multi-frame LiDAR points of dynamic objects and static scenes separately. In addition, we adopt Poisson Reconstruction [24] to fill up the holes and voxelize the obtained mesh to get dense volumetric occupancy. Since the LiDAR scans surface points, our method also generates surface occupancy.

### 4.1. Multi-frame Point Cloud Stitching

We propose a two-stream pipeline to stitch static scenes and objects separately and merge them into a complete scene before voxelization. Specifically, for each frame, we first cut out movable objects from the LiDAR points according to 3D bounding box labels, so that we can obtain the 3D points of a static scene and movable objects. After travers-

ing all frames in the scene, we integrate the collected static scene segments and object segments into a set respectively. To combine the multi-frame segments, we then transform their coordinates into the world coordinate system via the known calibrated matrices and ego-poses. We denote the transformed static scene segments and object segments as $P_{ss} = \{P_{ss}^1, P_{ss}^2, \cdots P_{ss}^n\}$ and $P_{os} = \{P_{os}^1, P_{os}^2, \cdots P_{os}^m\}$, where $n$ and $m$ are the numbers of frames and objects in the sequence, respectively. Note that the same objects in different frames can be recognized according to the bounding box index. Therefore, we can represent the whole static scene as $P_s = [P_{ss}^1, P_{ss}^2, \cdots P_{ss}^n]$ while the objects as $P_o = [P_{os}^1, P_{os}^2, \cdots P_{os}^m]$, respectively, where $[\cdot]$ is the concatenation operator. Finally, according to the objects' locations and ego-pose of the current frame, the 3D points of this frame can be obtained by merging static scene and objects: $P = [T_s(Ps), T_o(Po)]$, where $T_s$ and $T_o$ are the transformations of static scenes and objects from world coordinate system to the current frame coordinate system. In this way, the occupancy labels of the current frame leverage the LiDAR points of all frames in the sequence.

### 4.2. Densifying with Poisson Reconstruction

While the density of $P$ is much larger than a single-frame LiDAR, there still exists many interspaces and the points are not evenly distributed, which is caused by the limited LiDAR beams. To address this, we first compute the normal vectors according to the spatial distribution in local neighborhoods. Then we reconstruct $P$ to a triangular mesh $\mathcal{M}$ via Poisson Surface Reconstruction [24], whose input is point cloud with normal vectors, and the output is a triangular mesh (see Figure 4). The obtained mesh $\mathcal{M} = \{\mathcal{V}, \mathcal{E}\}$ fills up the holes of point clouds with evenly distributed vertices $\mathcal{V}$, so that we can further convert the mesh into dense voxel $V_d$.

### 4.3. Semantic Labeling with NN Algorithm

Having obtained the occupancy of dense voxels $V_d$, we aim to assign semantic labels to each voxel, since position reconstruction can only be applied to 3D space, not semantic space. To this end, we propose to leverage Nearest Neighbors (NN) algorithm to search the nearest semantic label for each voxel. Specifically, we first voxelize $P$ with semantics into voxels $V_s$, which are sparser than $V_d$ due to limited LiDAR beams. Then for each occupied voxel in $V_d$, we use NN to search the nearest voxel in $V_s$ and assign the semantic label to it. Note that this process can be accelerated by parallel computing on the GPU. Thus, all occupied voxels in $V_d$ can obtain their semantic labels from $V_s$.

Figure 5 shows a detailed visual comparison between single-frame LiDAR points, sparse occupancy labels and dense occupancy labels. We observe that our dense voxels can provide much more realistic occupancy labels with

| Method | SC IoU | SSC mIoU | barrier | bicycle | bus | car | const. veh. | motorcycle | pedestrian | traffic cone | trailer | truck | drive. suf. | other flat | sidewalk | terrain | manmade | vegetation |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MonoScene [8] | 23.96 | 7.31 | 4.03 | 0.35 | 8.00 | 8.04 | 2.90 | 0.28 | 1.16 | 0.67 | 4.01 | 4.35 | 27.72 | 5.20 | 15.13 | 11.29 | 9.03 | 14.86 |
| Atlas [37] | 28.66 | 15.00 | 10.64 | 5.68 | 19.66 | 24.94 | 8.90 | 8.84 | 6.47 | 3.28 | 10.42 | 16.21 | 34.86 | 15.46 | 21.89 | 20.95 | 11.21 | 20.54 |
| BEVFormer [29] | 30.50 | 16.75 | 14.22 | 6.58 | 23.46 | 28.28 | 8.66 | 10.77 | 6.64 | 4.05 | 11.20 | 17.78 | 37.28 | 18.00 | 22.88 | 22.17 | 13.80 | 22.21 |
| TPVFormer [22] | 11.51 | 11.66 | 16.14 | 7.17 | 22.63 | 17.13 | 8.83 | 11.39 | 10.46 | 8.23 | 9.43 | 17.02 | 8.07 | 13.64 | 13.85 | 10.34 | 4.90 | 7.37 |
| TPVFormer* | 30.86 | 17.10 | 15.96 | 5.31 | 23.86 | 27.32 | 9.79 | 8.74 | 7.09 | 5.20 | 10.97 | 19.22 | **38.87** | 21.25 | 24.26 | **23.15** | 11.73 | 20.81 |
| SurroundOcc | 31.49 | 20.30 | 20.59 | 11.68 | 28.06 | 30.86 | 10.70 | 15.14 | 14.09 | 12.06 | 14.38 | 22.26 | 37.29 | **23.70** | 24.49 | 22.77 | **14.89** | 21.86 |

Table 1. **3D semantic occupancy prediction results on nuScenes validation set.** Except TPVFormer [22], all methods are trained with dense occupancy labels. To fairly compare, we further use dense ground truth to train the TPVFormer, which is denoted as TPVFormer*.

| Acc | $\mathrm{mean}_{p \in P}(\min_{p^* \in P^*} ||p - p^*||)$ |
|---|---|
| Comp | $\mathrm{mean}_{p^* \in P^*}(\min_{p \in P} ||p - p^*||)$ |
| Prec | $\mathrm{mean}_{p \in P}(\min_{p^* \in P^*} ||p - p^*|| < 0.5)$ |
| Recal | $\mathrm{mean}_{p^* \in P^*}(\min_{p \in P} ||p - p^*|| < 0.5)$ |
| CD | Acc + Comp |
| F-score | $(2 \times \mathrm{Prec} \times \mathrm{Recal})/(\mathrm{Prec} + \mathrm{Recal})$ |

Table 2. Evaluation metrics for 3D scene reconstruction. $p$ and $p^*$ are the predicted and ground truth point clouds.

clear semantic boundaries.

We think it is not trivial to propagate the sparse semantic label to a dense one since it is a ill-posed problem. The proposed multi-frame point cloud stitching can aggregate multi-frame semantic information and provide dense enough reference points for NN. However, we find that NN is sensitive to the annotation noise in original LiDAR semantic labels, and we will try to solve it as the future work.

## 5. Experiments

### 5.1. Experimental Setup

**Dataset:** We conduct multi-camera experiments on nuScenes dataset [7], which is a large-scale autonomous driving dataset. Since the 3D semantic and 3D detection labels are unavailable in test set and we cannot generate dense occupancy labels, we use the training set to train the model and validation set for evaluation. The occupancy prediction range is set as $[-50m, 50m]$ for $X, Y$ axis and $[-5m, 3m]$ for $Z$ axis. The final output occupancy has the shape 200x200x16 with 0.5m voxel size. The input image resolution is 1600x900.

To further demonstrate the effectiveness of our method, we conduct monocular semantic scene completion experiment on SemanticKITTI dataset [2]. SemanticKITTI has annotated outdoor LiDAR scans with 21 semantic labels. The ground truth is voxelized as 256x256x32 grid with 0.2m voxel size. We evaluate our model on the test set.

**Implementation Details:** The whole network architecture obtains $M = 4$ levels and we do not add the skip connection in level 0. For nuScenes dataset, we adopt ResNet101-DCN [19, 14] with the initial weight from FCOS3D [53] as the backbone to extract image features. The features of stage 1,2,3 are fed to FPN [31] and used as multi-scale image features. The number of 2D-3D spatial attention layers are set as 1, 3, 6 for three levels. For SemanticKITTI dataset, following MonoScene [8], we use a pretrained EfficientNetB7 [50] as the backbone to generate multi-scale image features. We also adopt FPN to further fuse the features of different levels. We set the number of 2D-3D spatial attention layers as 1, 3, 8. All experiments are conducted on 8 RTX 3090s. For Possion Reconstruction, we accumulate both key-frame and non key-frame data. In details, we use LiDAR frames at 20Hz, resulting in 400 frames and around 13 million points for a 20 second sequence.

**Evaluation Metrics:** For 3D semantic occupancy prediction, we use the intersection over union (IoU) of occupied voxels, ignoring their semantic class as the evaluation metric of the scene completion (SC) task and the mIoU of all semantic classes for the SSC task.

$$\mathrm{IoU} = \frac{TP}{TP + FP + FN}$$
$$\mathrm{mIoU} = \frac{1}{C} \sum_{i=1}^{C} \frac{TP_i}{TP_i + FP_i + FN_i} \quad (4)$$

where $TP$, $FP$, $FN$ indicate the number of true positive, false positive, and false negative predictions. $C$ is the class number.

For 3D scene reconstruction, we first convert occupancy prediction to point clouds. Following [37, 6], we use 3D metrics for evaluation, which is shown in Table 9. Chamfer distance (CD) and F-score are the main metrics since they consider both precision and recall. Please refer to the supplementary for more evaluation metric details. For all evaluation, we adopt dense occupancy as ground truth since sparse LiDAR points cannot fully evaluate the quality of occupancy reconstruction.

| Method | SC IoU | SSC mIoU | road (15.30%) | sidewalk (11.1%) | parking (1.12%) | other-grnd (0.56%) | building (14.1%) | car (3.92%) | truck (0.16%) | bicycle (0.03%) | motorcycle (0.03%) | other-veh. (0.20%) | vegetation (39.3%) | trunk (0.51%) | terrain (9.17%) | person (0.07%) | bicyclist (0.07%) | motorcyclist (0.05%) | fence (3.90%) | pole (0.29%) | traf.-sign (0.08%) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| LMSCNet [46] | 31.38 | 7.07 | 46.70 | 19.50 | 13.50 | 3.10 | 10.30 | 14.30 | 0.30 | 0.00 | 0.00 | 0.00 | 10.80 | 0.00 | 10.40 | 0.00 | 0.00 | 0.00 | 5.40 | 0.00 | 0.00 |
| 3DSketch [11] | 26.85 | 6.23 | 37.70 | 19.80 | 0.00 | 0.00 | 12.10 | 17.10 | 0.00 | 0.00 | 0.00 | 0.00 | 12.10 | 0.00 | 16.10 | 0.00 | 0.00 | 0.00 | 3.40 | 0.00 | 0.00 |
| AICNet [27] | 23.93 | 7.09 | 39.30 | 18.30 | 19.80 | 1.60 | 9.60 | 15.30 | 0.70 | 0.00 | 0.00 | 0.00 | 9.60 | 1.90 | 13.50 | 0.00 | 0.00 | 0.00 | 5.00 | 0.10 | 0.00 |
| JS3C-Net [59] | 34.00 | 8.97 | 47.30 | 21.70 | 19.90 | 2.80 | 12.70 | 20.10 | 0.80 | 0.00 | 0.00 | 4.10 | 14.20 | 3.10 | 12.40 | 0.00 | 0.20 | 0.20 | 8.70 | 1.90 | 0.30 |
| MonoScene [8] | 34.16 | 11.08 | 54.70 | 27.10 | 24.80 | 5.70 | 14.40 | 18.80 | 3.30 | 0.50 | 0.70 | **4.40** | **14.90** | 2.40 | 19.50 | 1.00 | 1.40 | **0.40** | 11.10 | 3.30 | 2.10 |
| TPVFormer [22] | 34.25 | 11.26 | 55.10 | 27.20 | 27.40 | 6.50 | 14.80 | 19.20 | **3.70** | 1.00 | 0.50 | 2.30 | 13.90 | 2.60 | **20.40** | 1.10 | **2.40** | 0.30 | 11.00 | 2.90 | 1.50 |
| SurroundOcc | **34.72** | **11.86** | **56.90** | **28.30** | **30.20** | **6.80** | **15.20** | **20.60** | 1.40 | **1.60** | **1.20** | 4.40 | 14.90 | **3.40** | 19.30 | **1.40** | 2.00 | 0.10 | **11.30** | **3.90** | **2.40** |

Table 3. **Monocular Semantic scene completion results on SemanticKITTI test set.** For fair comparison, we use the performances of RGB-inferred versions of the first four methods, which are reported in MonoScene [8]. Although our method is not designed for monocular perception, we still outperform other methods for a large margin.
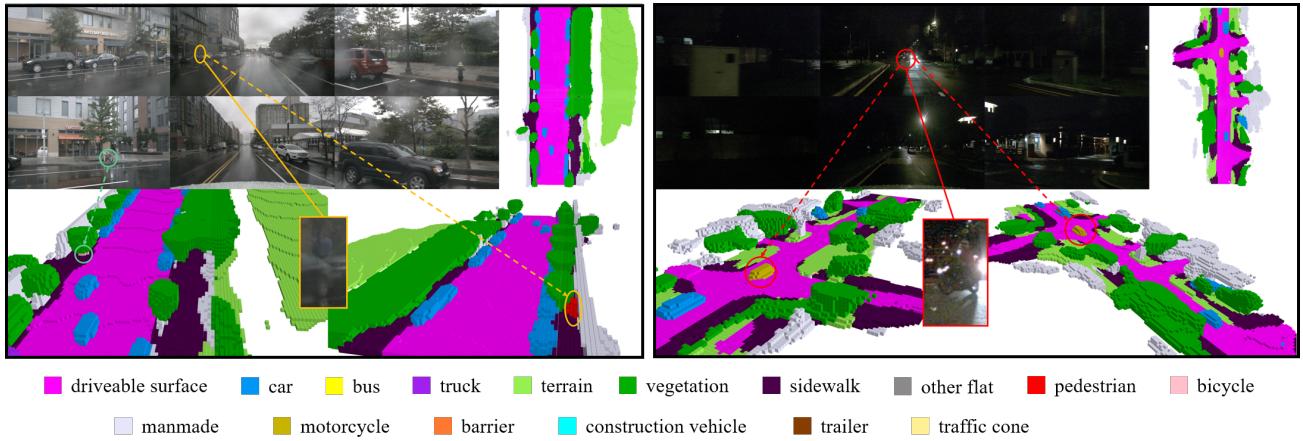


Figure 6. An example of challenging scenes. Although the quality of RGB images degrades in rainy days and nights, our method can still predict detailed occupancy. **Better viewed when zoomed in.**

| Method | Acc ↓ | Comp ↓ | Prec ↑ | Recall ↑ | CD ↓ | F-score ↑ |
|---|---|---|---|---|---|---|
| SurroundDepth [57] | 1.747 | 1.384 | 0.261 | 0.353 | 3.130 | 0.293 |
| AdaBins [4] | 1.989 | 1.287 | 0.233 | 0.347 | 3.275 | 0.271 |
| NeWCRFs [64] | 2.163 | 1.233 | 0.214 | 0.348 | 3.396 | 0.257 |
| Atlas [37] | **0.679** | 1.685 | 0.407 | 0.546 | 2.365 | 0.458 |
| TransformerFusion [6] | 0.771 | 1.434 | 0.375 | 0.591 | 2.205 | 0.453 |
| SurroundOcc | 0.724 | **1.226** | **0.414** | **0.602** | **1.950** | **0.483** |

Table 4. **3D scene reconstruction results on nuScenes validation set**. F-score and CD are the main metrics.

## 5.2. 3D Semantic Occupancy Prediction

We first conduct multi-camera 3D semantic occupancy prediction on nuScenes [7] dataset and compare with several state-of-the-art methods [8, 29, 22, 37]. For MonoScene [8], we concatenate the multi-camera occupancy predictions and voxelize them with the 0.5m voxel size, which is same as our setting. For BEVFormer [29], we add a 3D segmentation head to predict semantic occupancy. As shown in Table 1, our method achieves state-of-the-art performance. We also show some qualitative results in Figure 6 and Figure 7. See supplementary material for more video demos and qualitative comparisons.

Especially in Figure 6, we show rainy day and night visualization. Although the quality of RGB images degrades in these two challenging scenes, our method can still predict fine details. The LiDAR sensor suffers from rainy days and easily misses points. With multi-frame aggregation and Possion Reconstruction, we dramatically densify the labels and provide strong supervision, which is crucial for the challenging scenarios. Moreover, we conduct color jitter augmentation, which increases the robustness of brightness change. We also note that some parts cannot be observed by LiDAR sensor, such as the back side of the car. Surprisingly, our model can predict the complete shape according to the surrounding information.

To further demonstrate the superiority of our method, we also conduct monocular 3D semantic scene completion on SemanticKITTI dataset [2]. Table 3 shows the results. Although our method is not designed for monocular perception and cross-view attention will be ineffective for the monocular setting, our method still achieves state-of-the-art performance on this benchmark.

| Method | SC IoU | SSC mIoU |
|---|---|---|
| w/o spatial attention | 29.78 | 17.34 |
| BEV-based attention | 30.45 | 18.94 |
| Ours | **31.49** | **20.30** |

Table 5. The ablation study of 2D-3D spatial attention. "w/o spatial attention" indicates that we average all multi-camera features in a grid.

| Method | SC IoU | SSC mIoU |
|---|---|---|
| w/o multi-scale structure | 30.41 | 18.22 |
| w/o multi-scale supervision | 31.16 | 19.73 |
| Ours | **31.49** | **20.30** |

Table 6. The ablation study of multi-scale occupancy prediction. "w/o multi-scale structure" means that we do not add multi-scale skip connection.

| Supervision | SC IoU | SSC mIoU |
|---|---|---|
| sparse LiDAR points | 11.96 | 12.17 |
| sparse occupancy labels | 30.58 | 18.83 |
| dense occupancy labels | **31.49** | **20.30** |

Table 7. The ablation study of dense occupancy supervision. The model trained with our dense occupancy ground truth is much better than that trained with sparse LiDAR points.

## 5.3. 3D Scene Reconstruction

Another important application of 3D occupancy prediction is 3D scene reconstruction. Due to this reason, we further evaluate 3D reconstruction performance without using multi-class semantic labels. We do the comparison with state-of-the-art multi-camera depth estimation methods (SurroundDepth [57]), monocular depth estimation methods (AdaBins [4] and NeWCRFs [64]) and 3D reconstruction method (Atlas [37] and TransformerFusion [6]). For the self-supervised methods SurroundDepth [57], we use depth ground truth to supervise them. To evaluate depth estimation methods in 3D space, following [37, 6], we run TSDF fusion [13, 38] to fuse multi-camera depth as point clouds. Note that we fuse the multi-camera depth maps of the same timestamp but not the multi-frame depths. Thus, there is no need to specially deal with movable objects. As shown in Table 4, SurroundOcc achieves state-of-the-art performance on most metrics, and outperforms other methods by a large margin, which verifies the effectiveness of the proposed method.

## 5.4. Ablation Study

**2D-3D Spatial Attention:** Table 5 shows the ablation results for 2D-3D spatial attention. Without spatial attention, we directly average all multi-camera features in a grid. However, we find this straightforward fusion method performs worse than spatial attention. The potential reason is

that the contribution of each view is different for a 3D grid. Moreover, the ablation study shows that 3D-based cross-view attention is more effective than BEV-based cross-view attention since it can preserve 3D space information.

**Multi-scale Occupancy Prediction:** We conduct an ablation study on multi-scale structure and multi-scale supervision in Table 6. The experimental results show that these two multi-scale designs can boost the performance. The multi-scale skip connection can help the network learn low-level fine-grained features, which is important for detailed high-resolution occupancy prediction. Moreover, the 3D volume features in all levels will be enhanced by multi-scale supervision.

**Dense Occupancy Supervision:** The results in Table 7 demonstrate the importance of using dense occupancy as ground truth. Compared with sparse LiDAR points, the sparse occupancy labels fuse multi-frame points and can provide more powerful supervision. The possion reconstruction and NN algorithm can fill the holes and further densify the occupancy labels, which will boost model's performance. Figure 7 can better illustrate the effectiveness of dense supervision. We can see that the model trained with our occupancy labels can predict much denser occupancy than that trained with LiDAR points.

## 5.5. Model Efficiency

We compare the inference time and inference memory of different methods in Table 8. The experiments are conducted on one RTX 3090 with six multi-camera images, whose resolutions are 1600x900. We find that our method can achieve both high performance and efficiency. Compared with BEVFormer, our method slightly increases inference time and memory and we think the increased burden is acceptable

| Method | Latency (s) | Memory (G) |
|---|---|---|
| SurroundDepth [57] | 0.73 | 12.4 |
| NeWCRFs [64] | 1.07 | 14.5 |
| Adabins [4] | 0.75 | 15.5 |
| BEVFormer [29] | **0.31** | **4.5** |
| TPVFormer [22] | 0.32 | 5.1 |
| MonoScene [8] | 0.87 | 20.3 |
| Ours | 0.34 | 5.9 |

Table 8. The model efficiency of different methods. The experiments are conducted on one RTX 3090 with six multi-camera images, whose resolutions are 1600x900.

## 6. Limitations and Future Work

Currently, we only explore single-frame occupancy prediction. However, for the downstream modules, *e.g.* motion prediction and planning, occupancy flow is more important. As the future work, we will design a framework to build

FRONT LEFT  FRONT  FRONT RIGHT  BACK RIGHT  BACK  BACK LEFT

Sparse LiDAR points   Occupancy prediction (with LiDAR GT)   Occupancy prediction (with dense GT)   Dense occupancy labels
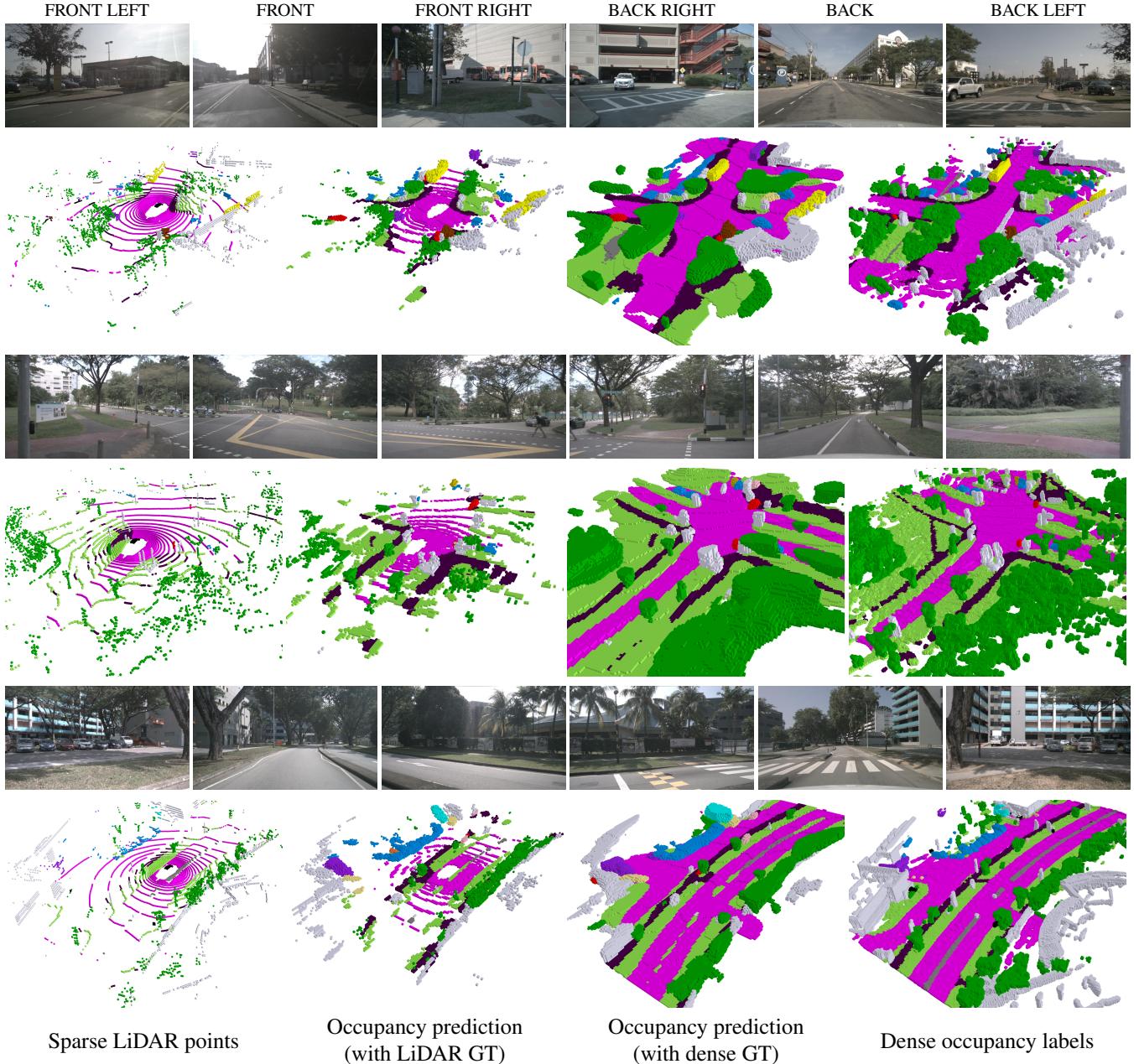
Figure 7. Visualizations on nuScenes validation set. Our generated dense occupancy labels are much denser than sparse LiDAR points. Trained with dense groundtruth, the network can predict better and denser occupancy. **Better viewed when zoomed in.**

occupancy flow dataset and utilize multi-frame surrounding images as the inputs. Moreover, LiDAR data is not always available. Self-supervised occupancy prediction with only RGB data is a valuable but challenging direction.

## 7. Conclusion

In this paper, we propose SurroundOcc for multi-camera 3D occupancy prediction. We utilize 2D-3D spatial attention to integrate 2D features to the 3D volume in a multi-scale fashion, which is further upsampled and fused by the 3D deconvolution layer. Moreover, we devise a pipeline to generate dense occupancy ground truth. We stitch multi-frame LiDAR points of dynamic objects and static scenes separately and utilize Poisson Reconstruction to fill the holes. The comparison on nuScenes and SemanticKITTI datasets demonstrates the superiority of our method.

## Acknowledgement

# Appendix

## A. Baseline Method Details

We compare with several baseline methods on nuScenes dataset, which can be roughly classified as four categories:

**Depth estimation:** SurroundDepth [57], AdaBins [4], NeWCRFs [64]. Since SurroundDepth method is multi-camera self-supervised method, we use depth groundtruth to supervise the network along with self-supervised photometric loss. AdaBins [4] and NeWCRFs [64] are the state-of-the-art depth estimation methods both in outdoor and indoor scenes. To implement these two methods, we use their official released code with the dataloader in SurroundDepth. The depth results are fused by the TSDF fusion algorithm [13, 38] with the voxel size $0.5m$, which is same to our method.

**3D scene reconstruction:** Atlas [37] and Transformerfusion [6]. These two methods are state-of-the-art indoor scene reconstruction methods. We use our dense occupancy groundtruth to supervise them instead of tsdf ground truth. To fairly compare, we also adopt ResNet101-DCN [19, 14] with the initial weight from FCOS3D [53] as the backbone to extract image features.

**Occupancy reconstruction:** MonoScene [8] and TPV-Former [22]. To extend MonoScene to multi-camera setting, we project occupancy labels to each camera's coordinate and the shape of each camera's prediction is $(128, 104, 16)$ with 0.5m voxel size. We fuse multi-camera results in LiDAR coordinate with camera extrinsics. The final result has the same shape and voxel size with ours. For TPVFormer, the resolution is set as 200x200x16 and the feature dimension is 64.

**BEV perception:** BEVFormer [29]. We use the full-resolution 200x200 BEV features. To lift BEV features to the 3D space, we split 256 dimensions BEV features as 16 grids and the feature of each grid has 16 dimensions. Then we adopt a 3D encoder-decoder network [37] as a segmentation head to predict occupancy. Following the setting in TPVFormer, we employ both cross entropy loss and lovasz-softmax [3] as the supervision signals.
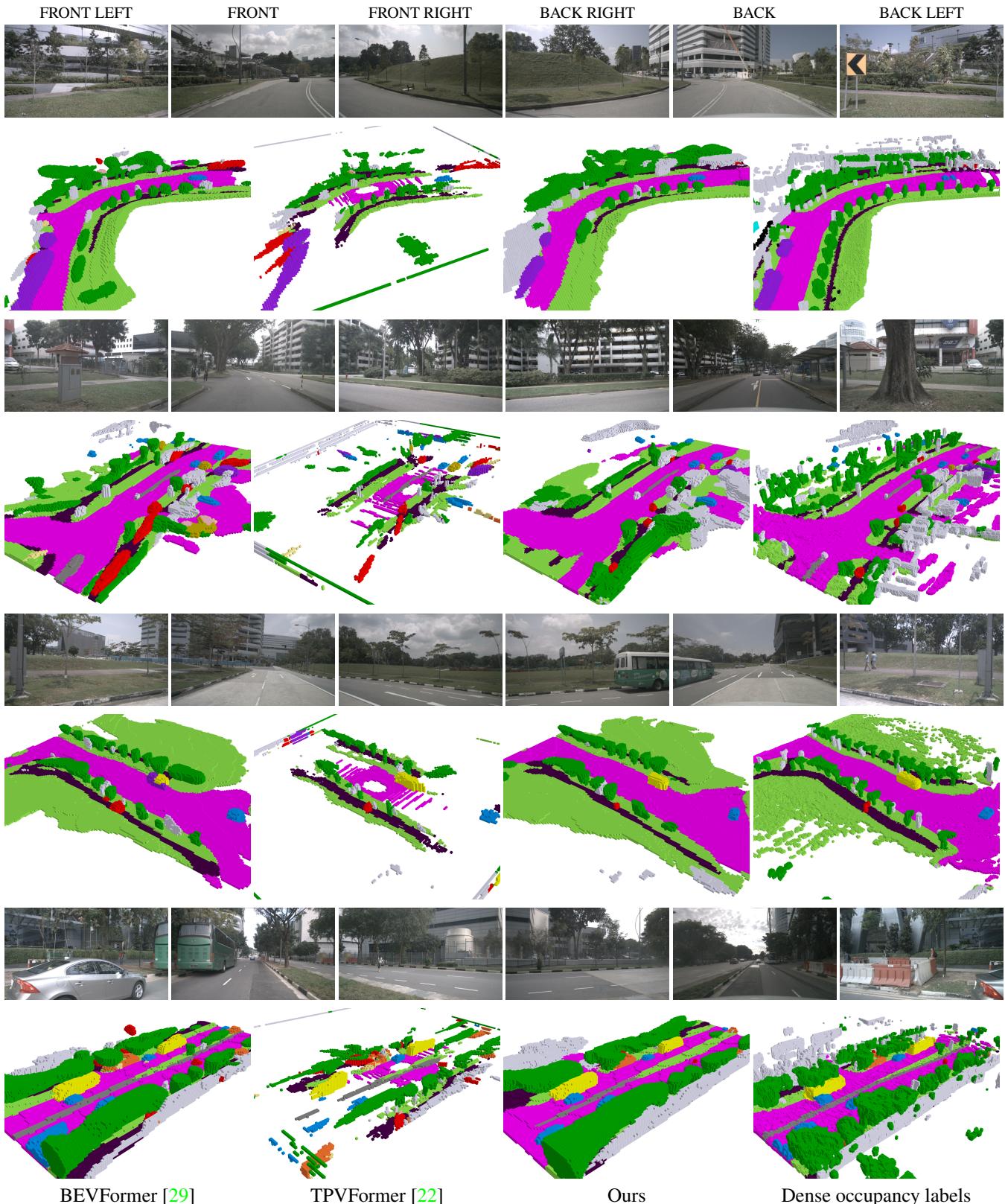
## B. More Visualizations

Figure 8 shows the qualitative comparison with other methods. We can see that our predictions are more accurate and denser. We also provide some video demos in the material. Specifically, 'demo-nuscenes' shows the results on nuScenes validation set and 'demo-gt' visualizes our generated groundtruth. 'demo-comparison' illustrates the comparison with other methods and 'demo-wild' shows the occupancy predictions on Beijing street (trained on nuScenes training set).

| | |
|---|---|
| Acc | $\text{mean}_{p \in P}(\min_{p^* \in P^*} ||p - p^*||)$ |
| Comp | $\text{mean}_{p^* \in P^*}(\min_{p \in P} ||p - p^*||)$ |
| Prec | $\text{mean}_{p \in P}(\min_{p^* \in P^*} ||p - p^*|| < 0.5)$ |
| Recal | $\text{mean}_{p^* \in P^*}(\min_{p \in P} ||p - p^*|| < 0.5)$ |
| CD | Acc + Comp |
| F-score | $(2 \times \text{Prec} \times \text{Recal})/(\text{Prec} + \text{Recal})$ |

Table 9. Evaluation metrics for 3D scene reconstruction. $p$ and $p^*$ are the predicted and ground truth point clouds.

## References

[1] Adil Kaan Akan and Fatma Güney. Stretchbev: Stretching future instance prediction spatially and temporally. *arXiv preprint arXiv:2203.13641*, 2022. 1

[2] Jens Behley, Martin Garbade, Andres Milioto, Jan Quenzel, Sven Behnke, Cyrill Stachniss, and Jurgen Gall. Semantickitti: A dataset for semantic scene understanding of lidar sequences. In *ICCV*, pages 9297–9307, 2019. 2, 5, 6, 7

[3] Maxim Berman, Amal Rannen Triki, and Matthew B Blaschko. The lovász-softmax loss: A tractable surrogate for the optimization of the intersection-over-union measure in neural networks. In *CVPR*, pages 4413–4421, 2018. 10

[4] Shariq Farooq Bhat, Ibraheem Alhashim, and Peter Wonka. Adabins: Depth estimation using adaptive bins. In *CVPR*, pages 4009–4018, 2021. 7, 8, 10

[5] Jia-Wang Bian, Zhichao Li, Naiyan Wang, Huangying Zhan, Chunhua Shen, Ming-Ming Cheng, and Ian Reid. Unsupervised Scale-consistent Depth and Ego-motion Learning from Monocular Video. In *NeurIPS*, pages 35–45, 2019. 2

[6] Aljaz Bozic, Pablo Palafox, Justus Thies, Angela Dai, and Matthias Nießner. Transformerfusion: Monocular rgb scene reconstruction using transformers. *Advances in Neural Information Processing Systems*, 34:1403–1414, 2021. 2, 6, 7, 8, 10

[7] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *CVPR*, pages 11621–11631, 2020. 2, 6, 7

[8] Anh-Quan Cao and Raoul de Charette. Monoscene: Monocular 3d semantic scene completion. In *CVPR*, pages 3991–4001, 2022. 2, 3, 4, 6, 7, 8, 10

[9] Jia-Ren Chang and Yong-Sheng Chen. Pyramid stereo matching network. In *CVPR*, pages 5410–5418, 2018. 2

[10] Xiaokang Chen, Kwan-Yee Lin, Chen Qian, Gang Zeng, and Hongsheng Li. 3d sketch-aware semantic scene completion via semi-supervised structure prior. In *CVPR*, pages 4193–4202, 2020. 2

[11] Xiaokang Chen, Kwan-Yee Lin, Chen Qian, Gang Zeng, and Hongsheng Li. 3d sketch-aware semantic scene completion via semi-supervised structure prior. In *CVPR*, pages 4193–4202, 2020. 7

[12] Ran Cheng, Ryan Razani, Ehsan Taghavi, Enxu Li, and Bingbing Liu. 2-s3net: Attentive feature fusion with adaptive feature selection for sparse semantic segmentation network. In *CVPR*, pages 12547–12556, 2021. 2

| FRONT LEFT | FRONT | FRONT RIGHT | BACK RIGHT | BACK | BACK LEFT |

| BEVFormer [29] | TPVFormer [22] | Ours | Dense occupancy labels |

Figure 8. Qualitaative comparison on nuScenes validation set. Our mrthod can predict more accurate and denser occupancy. **Better viewed when zoomed in.**

[13] Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *SIGGRAPH*, pages 303–312, 1996. 8, 10

[14] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *ICCV*, pages 764–773, 2017. 6, 10

[15] David Eigen and Rob Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *ICCV*, pages 2650–2658, 2015. 2

[16] Lue Fan, Xuan Xiong, Feng Wang, Naiyan Wang, and Zhaoxiang Zhang. Rangedet: In defense of range view for lidar-based 3d object detection. In *ICCV*, pages 2918–2927, 2021. 2

[17] Clément Godard, Oisin Mac Aodha, Michael Firman, and Gabriel J Brostow. Digging into self-supervised monocular depth estimation. In *ICCV*, pages 3828–3838, 2019. 2

[18] Vitor Guizilini, Igor Vasiljevic, Rares Ambrus, Greg Shakhnarovich, and Adrien Gaidon. Full surround monodepth from multiple cameras. *IEEE Robotics and Automation Letters*, 7(2):5397–5404, 2022. 1

[19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 3, 6, 10

[20] Anthony Hu, Zak Murez, Nikhil Mohan, Sofía Dudas, Jeffrey Hawke, Vijay Badrinarayanan, Roberto Cipolla, and Alex Kendall. Fiery: Future instance prediction in bird's-eye view from surround monocular cameras. In *ICCV*, pages 15273–15282, 2021. 1

[21] Junjie Huang, Guan Huang, Zheng Zhu, and Dalong Du. Bevdet: High-performance multi-camera 3d object detection in bird-eye-view. *arXiv preprint arXiv:2112.11790*, 2021. 1, 2

[22] 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. 2, 6, 7, 8, 10, 11

[23] Mengqi Ji, Juergen Gall, Haitian Zheng, Yebin Liu, and Lu Fang. Surfacenet: An end-to-end 3d neural network for multiview stereopsis. In *ICCV*, pages 2307–2315, 2017. 2

[24] Michael Kazhdan, Matthew Bolitho, and Hugues Hoppe. Poisson surface reconstruction. In *Proceedings of the fourth Eurographics symposium on Geometry processing*, volume 7, 2006. 2, 5

[25] Jin Han Lee, Myung-Kyu Han, Dong Wook Ko, and Il Hong Suh. From big to small: Multi-scale local planar guidance for monocular depth estimation. *arXiv preprint arXiv:1907.10326*, 2019. 2

[26] Jie Li, Kai Han, Peng Wang, Yu Liu, and Xia Yuan. Anisotropic convolutional networks for 3d semantic scene completion. In *CVPR*, pages 3351–3359, 2020. 2

[27] Jie Li, Kai Han, Peng Wang, Yu Liu, and Xia Yuan. Anisotropic convolutional networks for 3d semantic scene completion. In *CVPR*, pages 3351–3359, 2020. 7

[28] Yinhao 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. 1, 2

[29] 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. In *ECCV*, 2022. 1, 2, 4, 6, 7, 8, 10, 11

[30] Tingting Liang, Hongwei Xie, Kaicheng Yu, Zhongyu Xia, Zhiwei Lin, Yongtao Wang, Tao Tang, Bing Wang, and Zhi Tang. Bevfusion: A simple and robust lidar-camera fusion framework. *arXiv preprint arXiv:2205.13790*, 2022. 2

[31] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, pages 2117–2125, 2017. 6

[32] Venice Erin Liong, Thi Ngoc Tho Nguyen, Sergi Widjaja, Dhananjai Sharma, and Zhuang Jie Chong. Amvnet: Assertion-based multi-view fusion network for lidar semantic segmentation. *arXiv preprint arXiv:2012.04934*, 2020. 2

[33] Yingfei Liu, Tiancai Wang, Xiangyu Zhang, and Jian Sun. Petr: Position embedding transformation for multi-view 3d object detection. In *ECCV*, pages 531–548. Springer, 2022. 1

[34] 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

[35] Zhijian Liu, Haotian Tang, Alexander Amini, Xinyu Yang, Huizi Mao, Daniela Rus, and Song Han. Bevfusion: Multi-task multi-sensor fusion with unified bird's-eye view representation. *arXiv preprint arXiv:2205.13542*, 2022. 1, 2

[36] Xinzhu Ma, Zhihui Wang, Haojie Li, Pengbo Zhang, Wanli Ouyang, and Xin Fan. Accurate monocular 3d object detection via color-embedded 3d reconstruction for autonomous driving. In *ICCV*, 2019. 2

[37] Zak Murez, Tarrence van As, James Bartolozzi, Ayan Sinha, Vijay Badrinarayanan, and Andrew Rabinovich. Atlas: End-to-end 3d scene reconstruction from posed images. In *ECCV*, pages 414–431, 2020. 2, 3, 6, 7, 8, 10

[38] Richard A Newcombe, Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J Davison, Pushmeet Kohi, Jamie Shotton, Steve Hodges, and Andrew Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. In *ISMAR*, pages 127–136, 2011. 8, 10

[39] Dennis Park, Rares Ambrus, Vitor Guizilini, Jie Li, and Adrien Gaidon. Is pseudo-lidar needed for monocular 3d object detection? In *ICCV*, 2021. 2

[40] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *CVPR*, pages 165–174, 2019. 2

[41] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *ECCV*, pages 523–540. Springer, 2020. 2

[42] Jonah Philion and Sanja Fidler. Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d. In *ECCV*, pages 194–210, 2020. 2

[43] Anurag Ranjan, Varun Jampani, Lukas Balles, Kihwan Kim, Deqing Sun, Jonas Wulff, and Michael J Black. Competitive Collaboration: Joint Unsupervised Learning of Depth,

Camera Motion, Optical Flow and Motion Segmentation. In *CVPR*, pages 12240–12249, 2019. 2

[44] Cody Reading, Ali Harakeh, Julia Chae, and Steven L Waslander. Categorical depth distribution network for monocular 3d object detection. In *CVPR*, 2021. 2

[45] 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. 2

[46] Luis Roldao, Raoul de Charette, and Anne Verroust-Blondet. Lmscnet: Lightweight multiscale 3d semantic completion. In *3DV*, pages 111–119. IEEE, 2020. 7

[47] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *ICCV*, pages 2304–2314, 2019. 2

[48] Shuran Song, Fisher Yu, Andy Zeng, Angel X Chang, Manolis Savva, and Thomas Funkhouser. Semantic scene completion from a single depth image. In *CVPR*, pages 1746–1754, 2017. 2

[49] Jiaming Sun, Yiming Xie, Linghao Chen, Xiaowei Zhou, and Hujun Bao. Neuralrecon: Real-time coherent 3d reconstruction from monocular video. In *CVPR*, pages 15598–15607, 2021. 2

[50] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, pages 6105–6114. PMLR, 2019. 6

[51] Haotian Tang, Zhijian Liu, Shengyu Zhao, Yujun Lin, Ji Lin, Hanrui Wang, and Song Han. Searching efficient 3d architectures with sparse point-voxel convolution. In *ECCV*, pages 685–702, 2020. 2

[52] Tai Wang, ZHU Xinge, Jiangmiao Pang, and Dahua Lin. Probabilistic and geometric depth: Detecting objects in perspective. In *CoRL*, 2022. 2

[53] Tai Wang, Xinge Zhu, Jiangmiao Pang, and Dahua Lin. Fcos3d: Fully convolutional one-stage monocular 3d object detection. In *ICCV*, 2021. 2, 6, 10

[54] Yan Wang, Wei-Lun Chao, Divyansh Garg, Bharath Hariharan, Mark Campbell, and Kilian Q Weinberger. Pseudo-lidar from visual depth estimation: Bridging the gap in 3d object detection for autonomous driving. In *CVPR*, 2019. 2

[55] Yue Wang, Vitor Guizilini, Tianyuan Zhang, Yilun Wang, Hang Zhao, and Justin M. Solomon. Detr3d: 3d object detection from multi-view images via 3d-to-2d queries. In *CoRL*, 2021. 2

[56] Zhongdao Wang, Luming Tang, Xihui Liu, Zhuliang Yao, Shuai Yi, Jing Shao, Junjie Yan, Shengjin Wang, Hongsheng Li, and Xiaogang Wang. Orientation invariant feature embedding and spatial temporal regularization for vehicle re-identification. In *ICCV*, pages 379–387, 2017. 2

[57] Yi Wei, Linqing Zhao, Wenzhao Zheng, Zheng Zhu, Yongming Rao, Guan Huang, Jiwen Lu, and Jie Zhou. Surround-depth: Entangling surrounding views for self-supervised multi-camera depth estimation. In *CoRL*, 2022. 1, 2, 7, 8, 10

[58] 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

[59] 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 *AAAI*, volume 35, pages 3101–3109, 2021. 7

[60] 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

[61] Maosheng Ye, Rui Wan, Shuangjie Xu, Tongyi Cao, and Qifeng Chen. Drinet++: Efficient voxel-as-point point cloud segmentation. *arXiv preprint arXiv: 2111.08318*, 2021. 2

[62] Zhichao Yin and Jianping Shi. Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In *CVPR*, pages 1983–1992, 2018. 2

[63] Chen Yongjian, Tai Lei, Sun Kai, and Li Mingyang. Monopair: Monocular 3d object detection using pairwise spatial relationships. In *CVPR*, 2020. 2

[64] Weihao Yuan, Xiaodong Gu, Zuozhuo Dai, Siyu Zhu, and Ping Tan. New crfs: Neural window fully-connected crfs for monocular depth estimation. In *CVPR*, 2022. 7, 8, 10

[65] Yunpeng Zhang, Zheng Zhu, Wenzhao Zheng, Junjie Huang, Guan Huang, Jie Zhou, and Jiwen Lu. Beverse: Unified perception and prediction in birds-eye-view for vision-centric autonomous driving. *arXiv preprint arXiv:2205.09743*, 2022. 1, 2

[66] Zhenyu Zhang, Chunyan Xu, Jian Yang, Junbin Gao, and Zhen Cui. Progressive hard-mining network for monocular depth estimation. *TIP*, 27(8):3691–3702, 2018. 2

[67] Wang Zhao, Shaohui Liu, Yezhi Shu, and Yong-Jin Liu. Towards Better Generalization: Joint Depth-Pose Learning without PoseNet. In *CVPR*, 2020. 2

[68] 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. 2

[69] Junsheng Zhou, Yuwang Wang, Kaihuai Qin, and Wenjun Zeng. Moving Indoor: Unsupervised Video Depth Learning in Challenging Environments. In *ICCV*, pages 8618–8627, 2019. 2

[70] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *CVPR*, pages 1851–1858, 2017. 2

[71] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. In *ICLR*, 2020. 4

[72] Xinge Zhu, Hui Zhou, Tai Wang, Fangzhou Hong, Yuexin Ma, Wei Li, Hongsheng Li, and Dahua Lin. Cylindrical and asymmetrical 3d convolution networks for lidar segmentation. In *CVPR*, pages 9939–9948, 2021. 2