# OctreeOcc: Efficient and Multi-Granularity Occupancy Prediction Using Octree Queries

**Yuhang Lu[1], Xinge Zhu[3], Tai Wang[2,†], Yuexin Ma[1,†]**

[1] ShanghaiTech University  [2] Shanghai AI Laboratory  [3] The Chinese University of Hong Kong

{luyh2,mayuexin}@shanghaitech.edu.cn, taiwang.me@gmail.com

## Abstract

*Occupancy prediction has increasingly garnered attention in recent years for its fine-grained understanding of 3D scenes. Traditional approaches typically rely on dense, regular grid representations, which often leads to excessive computational demands and a loss of spatial details for small objects. This paper introduces OctreeOcc, an innovative 3D occupancy prediction framework that leverages the octree representation to adaptively capture valuable information in 3D, offering variable granularity to accommodate object shapes and semantic regions of varying sizes and complexities. In particular, we incorporate image semantic information to improve the accuracy of initial octree structures and design an effective rectification mechanism to refine the octree structure iteratively. Our extensive evaluations show that OctreeOcc not only surpasses state-of-the-art methods in occupancy prediction, but also achieves a $15\% - 24\%$ reduction in computational overhead compared to dense-grid-based methods.*

## 1. Introduction

Holistic 3D scene understanding is a pivotal aspect of a stable and reliable visual perception system, especially for real-world applications such as autonomous driving. Occupancy, as a classical representation, has been renascent recently with more datasets support and exploration in learning-based approaches. Such occupancy prediction tasks aim at partitioning the 3D scene into grid cells and predicting semantic labels for each voxel. It is particularly an essential solution for recognizing irregularly shaped objects and also enables the open-set understanding [34], further benefiting downstream tasks, like prediction and planning.

Existing occupancy prediction methods [8, 10, 19, 42, 45] typically construct dense grid representations, same as the occupancy ground truth. It is a natural and straightforward solution but does not take the property of 3D scene



(a) Sparsity of 3D scene    (b) Reference of octree in 2D

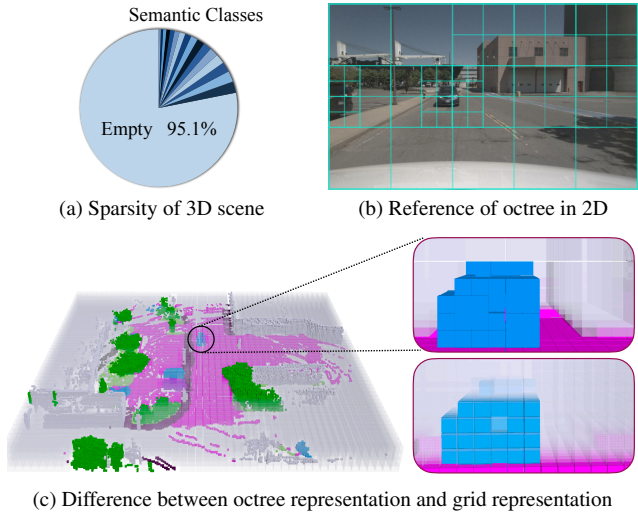(c) Difference between octree representation and grid representation

Figure 1. **Illustration of the sparsity of 3D scene and octree representation.** (a) shows that the semantic categories only account for a small fraction of the space, and that modeling empty regions densely affects computational efficiency. (b) and (c) demonstrate the superiority of octree representations, where we can apply different granularities for different scales of objects or spaces, which reduces computational overhead while preserving spatial information. Translucent areas represent empty voxels.

geometry into consideration. The 3D world is composed of foreground objects and background regions with shapes and sizes varying significantly. For example, in the occupancy ground truth, there can be more than 90% empty regions as shown in Fig. 1a. The regions occupied by objects such as buses and construction vehicles are also much larger than small objects like traffic cones. Thus, using the same resolution of voxels to represent the scene is sub-optimal. As a result, the aforementioned works struggle with large computational overhead. Other recent works attempted to mitigate the heavy memory footprint with other representations, such as 2D planes and coarse voxels exploited in TPV-Former [9] and PanoOcc [41]. VoxFormer [17] leverages the depth estimation to model non-empty regions. However,

---

they still suffer from the loss of spatial information (*i.e.* 2D planes) and rely on the depth estimation quality.

In contrast to previous attempts, we tackle this issue with the octree [25] representation that exactly accommodates object shapes and semantic regions with varying sizes. As a tree-based data structure, it recursively divides the 3D space into eight octants, thus allowing coarse spatial partition for large regions but fine-grained processing for small objects or complex details (Fig. 1b and 1c). Incorporating this representation into the occupancy prediction framework, we propose *OctreeOcc*, an efficient and multi-granularity framework using octree queries. It constructs octree queries by predicting the octree structure from the stored features of leaf nodes at each level of the tree. Such 3D structure prediction from 2D images is challenging due to the lack of depth and occlusion issues. To address this problem, we first propose **Semantic-guided Octree Initialization** that incorporates image semantic information to produce more accurate initial structures. To further improve the octree prediction, we devise an **Iterative Structure Rectification** module that predicts new octree structure from the encoded query to rectify the low-confidence region of the original prediction.

Our extensive evaluations against state-of-the-art occupancy prediction methods demonstrate that OctreeOcc outperforms others on nuScenes and SemanticKITTI datasets, reducing computational overhead by $15\% - 24\%$ for dense-grid representation-based methods. Ablation studies further validate the effectiveness of each module within our method. In summary, our contributions can be summarized as follows:

- We introduce OctreeOcc, a 3D occupancy prediction framework based on multi-granularity octree query. The scheme achieves spatial sparsification, and effectively reduces the number of voxels required to represent the scene while preserving the essential spatial information.
- We develop the semantic-guided octree initialization module and iterative structure rectification module to enable the network to have a good initialization and dynamically adjust the octree for a more effective representation.
- Comprehensive experiments demonstrate that OctreeOcc achieves state-of-the-art performance and reduces computational overhead, highlighting the feasibility and potential of octree structures in 3D occupancy prediction.

## 2. Related Work

### 2.1. Camera-based 3D Perception

Camera-based 3D perception has gained significant traction in recent years due to its ease of deployment, cost-effectiveness, and the preservation of intricate visual attributes. According to the view transformation paradigm, these methods can be categorized into three distinct types.

LSS-based approaches[8, 15, 16, 24, 28, 29] explicitly lift multi-view image features into 3D space through depth prediction. Another category of works[11, 18, 44] implicitly derives depth information by querying from 3D to 2D. Notably, projection-free methods[21–23, 39] have recently demonstrated exceptional performance. While commendable progress has been made in detection, this approach compromises the comprehensive representation of the overall scene structure within 3D space and proves less effective in recognizing irregularly shaped objects. Consequently, there is a burgeoning interest in methodologies aimed at acquiring a dense voxel representation through the camera, facilitating a more comprehensive understanding of 3D space.

### 2.2. 3D Occupancy Prediction

3D occupancy prediction involves the prediction of both occupancy and semantic attributes for all voxels encompassed within a three-dimensional scene, particularly valuable for autonomous vehicular navigation. Recently, some valuable datasets [33, 35, 40] have been proposed, boosting more and more research works [27, 34] in this field. Most of the research in occupancy prediction focuses on dense voxel modeling. MonoScene[3] pioneers a camera-based approach using a 3D UNet architecture. OccDepth[26] improves 2D-to-3D geometric projection using stereo-depth information. OccFormer[45] decomposes the 3D processing into the local and global transformer pathways along the horizontal plane. SurroundOcc[42] achieves fine-grained results with multiscale supervision. Symphonies[10] introduces instance queries to enhance scene representation.

Nevertheless, owing to the high resolution of regular voxel representation and sparse context distribution in 3D scenes, these methods encounter substantial computational overhead and efficiency issues. Some approaches recognize this problem and attempt to address it by reducing the number of modeled voxels. For instance, TPVFormer[9] employs a strategy of modeling the three-view 2D planes and subsequently recovering 3D spatial information from them. However, its performance degrades due to the lack of 3D information. PanoOcc[41] initially represents scenes at the coarse-grained level and then upsamples them to the fine-grained level, but the lack of information from coarse-grained modeling cannot be adequately addressed by the up-sampling process. VoxFormer[17] mitigates computational complexity by initially identifying non-empty regions through depth estimation and modeling only those specific areas. However, the effectiveness of this process is heavily contingent on the accuracy of depth estimation.

In contrast, our approach provides different granularity of modeling for different regions by predicting the octree structure, which reduces the number of voxels to be modeled while preserving the spatial information, thereby reducing the computational overhead and maintaining the pre-

diction accuracy.

## 2.3. Octree-Based 3D Representation

The octree technique[25] involves recursively partitioning three-dimensional space into eight octants. Widely utilized in computer graphics, it is integral to applications such as rendering, shape reconstruction, and collision detection. Due to its spatial efficiency and ease of implementation on GPU architectures, researchers have applied octrees in 3D point cloud learning and various related works. Some research[13, 37, 38] leverages the octree structure to efficiently represent and analyze the shape of the point cloud. OctFormer[36] and OcTr [46] utilize octree multi-granularity features to capture a comprehensive global context, thereby enhancing the efficiency of understanding point clouds at the scene level. Furthermore, certain studies [6, 30] adopt octree representation for effectively compressing point cloud data. These instances highlight the versatility and effectiveness of octree-based methodologies in point cloud analysis and processing applications. However, unlike constructing an octree from a 3D point cloud, predicting the octree structure of a 3D scene from images is a challenging task owing to the absence of explicit spatial information inherent in 2D representations. In order to address this complexity, we take the semantic information as a precursor to initializing the octree structure. Subsequently, we iteratively update the structure, ensuring a continual rectification that aligns more accurately with the scene.

## 3. Methods

### 3.1. Problem Steup

Camera-based occupancy prediction aims to predict the present occupancy state and semantics of each voxel grid within the scene using input from multi-view camera images. Specifically, we consider a set of $N$ multi-view images $I = \{I_i \in \mathbb{R}^{3 \times H \times W}\}_{i=1}^N$, together with camera intrinsics $K = \{K_i \in \mathbb{R}^{3 \times 3}\}_{i=1}^N$ and extrinsics $T = \{T_i \in \mathbb{R}^{4 \times 4}\}_{i=1}^N$ as input and the objective of the model is to predict the 3D occupancy $O \in \mathbb{R}^{X \times Y \times Z \times K}$, where $X$, $Y$, $Z$ denote the volume resolution, and $K$ represents the number of semantic categories.

### 3.2. Overall Architecture

As illustrated in Figure 2, OctreeOcc primarily comprises four modules. Initially, the image backbone extracts multi-scale features from multi-view images. Then we randomly initialize the dense occupancy queries $Q_{dense} \in \mathbb{R}^{X \times Y \times Z \times C}$, where $X$, $Y$, $Z$ are the spatial shape of occupancy queries, and C is the channel dimension. Subsequently, leveraging an image segmentation prior (Section 3.4), the octree initialization module transforms $Q_{dense}$ into sparse octree queries $Q_{octree}$ (Section 3.3). Facilitated by

the octree encoder, image features are conveyed to the octree query, and the octree structure undergoes continuous refinement (Section 3.5). Finally, an octree decoder translates octree queries back into dense grid queries, culminating in the generation of outputs for the occupancy prediction task.

### 3.3. Octree Query

Given that objects within a 3D scene exhibit varying granularities, employing dense queries for processing coarse-grained objects leads to performance wastage and inefficiency. To address this issue, we introduce sparse and multi-granularity octree queries, leveraging the octree structure. This approach enables the creation of flexible voxel representations tailored to semantic regions across different scales. To construct an octree query, we initiate the process by predicting the octree mask. This mask serves as a critical element in acquiring the octree structure, facilitating the conversion of $Q_{dense}$ into $Q_{octree}$. Denoted as $M_o \in \mathbb{R}^{\frac{X}{2^l}, \frac{Y}{2^l}, \frac{Z}{2^l}}$, where $l = 1, 2, \ldots, L - 1$, the octree mask signifies the likelihood that each voxel at every level undergoes a split into octants. Here, $L$ represents the depth of the octree.

Based on $M_o$ and a query selection ratio denoted as $\alpha = \alpha^1, \alpha^2, \ldots, \alpha^{L-1}$, which governs the number of subdivisions at each level, we can define the hierarchical structure of the octree, specifically determining how queries are assigned to respective levels. The procedure is as follows. For level 1, queries with $M_o^1$ values within the top $\alpha^1$ percentile are identified as candidates for octant subdivision. Subsequently, within octants that have undergone one previous subdivision at the next level, queries are once again selected based on their values falling within the top $\alpha^2$ percentile, initiating another round of subdivision. This process continues iteratively until reaching the final level of the octree.

By leveraging structural information, we can transform dense query into sparse and multi-granularity octree query. Specifically, we downsample the $Q_{dense}$ at each level using average pooling. Within each level, we selectively retain queries that have no need for further subdivision into octants (i.e., leaf queries). By storing all the constructed leaf queries, we form the $Q_{octree} \in \mathbb{R}^{N \times C}$, where $N = N_1 + N_2 + \ldots + N_L$ denotes the number of leaf queries. This information is updated in the octree encoder module along with the octree mask.

At the same time, leveraging the octree structure allows for the straightforward conversion of sparse octree queries into dense queries to match the desired output shape. Specifically, for a query $q_{octree}$ at level $l$ with the index $(a, b, c)$, its corresponding $8^{(L - l)}$ children nodes' indexes in level $L$ are given by $(a \times 2^{L-l} + a_{offset}, b \times 2^{L-l} + b_{offset}, c \times 2^{L-l} + c_{offset})$, where $a_{offset}, b_{offset}, c_{offset} = 0, 1, \ldots, 2^{L-l}$ and L denotes the
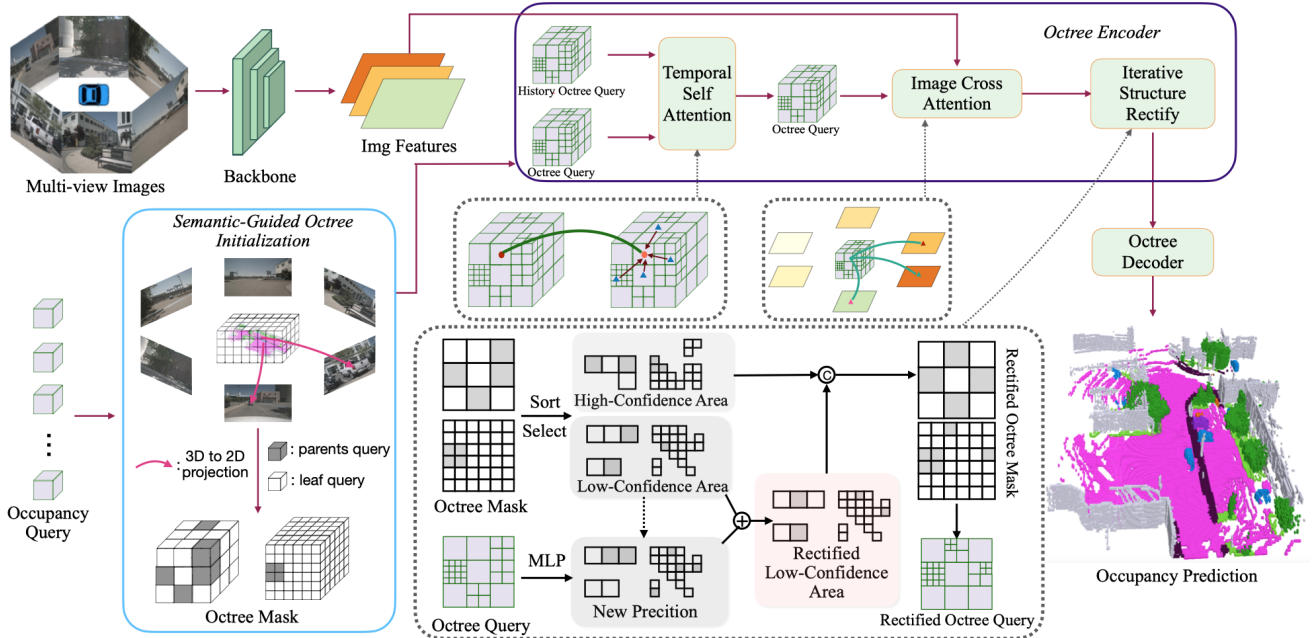
Figure 2. **Overall framework of OctreeOcc.** Given multi-view images, we extract multi-scale image features utilizing an image backbone. Subsequently, the initial octree structure is derived through image segmentation priors, and the transformation of dense queries into octree queries is effected. Following this, we concomitantly refine octree queries and rectify the octree structure through the octree encoder. Finally, we decode from the octree query and obtain occupancy prediction outcomes for this frame. For better visualisation, the diagram of Iterative Structure Rectification module shows octree query and mask in 2D form(quadtree).

depth of the octree. In this process, we assign $q_{octree}$'s feature to all of these positions. By repeating this procedure for all queries at each level, we effectively convert the $Q_{octree}$ into $Q_{dense}$.

### 3.4. Semantic-Guided Octree Initialization

Predicting octree structure from an initialised query via neural network can yield sub-optimal results due to the inherent lack of meaningful information in the query. To address this limitation, we employ the image prior as a pivotal reference point. Specifically, we adopt UNet[32] to segment the input multi-view images $I$ and obtain the segment map $I_{seg} \in \mathbb{R}^{H \times W \times K_{seg}}$. We then generate sampling points $\{p_i = (x_i, y_i, z_i), i \in X \times Y \times Z\}$ corresponding to a single occupancy query $Q^{p_i}$ at grid region center $p_i = (x_i, y_i, z_i)$, and then project these points to different image views. The projection between sampling points $p_i$ and its corresponding 2D reference point $(u_{ij}, v_{ij})$ on the $j$-th image view is formulated as:

$$\pi_j(p_i) = (u_{ij}, v_{ij}) \tag{1}$$

where $\pi_j(p_i)$ denotes the projection of the $i$-th sampling point at location $p_i$ on the $j$-th camera view.

We project the point $p_i$ onto the acquired semantic segmentation map $I_{seg}$ through the described projection process. To address the impact of the long-tailed distribution

inherent in the dataset, we assign varying weights to voxels projected onto different semantic categories. After obtaining the initial confidence for each voxel, we downsample it to different octree levels to formulate the initial octree mask. This mask will be further refined in the octree encoder.

### 3.5. Octree Encoder

Given octree queries $Q_{octree}$ and extracted image features $\mathbb{F}$, the octree encoder updates both the octree query features and the octree structure.

**Spatial-Tenporal Attention.** Referring to the querying paradigm in dense query-based methods[33, 41], we adopt efficient deformable attention[48] for image cross-attention(ICA) and temporal self-attention(TSA).

ICA is devised to enhance the interaction between multi-scale image features and octree queries. Specifically, for an octree query $q$, we can obtain its centre's 3D coordinate $(x, y, z)$ as reference point $Ref_{x,y,z}$. Then we project the 3d point to images like formula 1 and perform deformable attention:

$$ICA(q, \mathbf{F}) = \frac{1}{N} \sum_{n \in N} \sum_{m=1}^{M_1} DA(q, \pi_n(Ref_{x,y,z}), \mathbf{F}_n) \tag{2}$$

where $N$ denotes the camera view, $m$ indexes the reference points , and $M_1$ is the total number of sampling points for

each query. $\mathbf{F}_n$ is the image features of the $n$-th camera view. DA represents deformable attention.

In accurately representing the driving scene, temporal information plays a crucial role. Given the history octree queries $Q_{t-1}$, we align it to the current octree queries $Q_t$ by the ego-vehicle motion. To mitigate computational costs, we employ similar operation to ICA. Given octree queries $Q_t \in \mathbb{R}^{N,C}$, a single query with feature $q \in \mathbb{R}^C$ and the cooresponding coordinate $p = (x, y, z)$, the TSA is represented by:

$$TSA(q, p, Q_t) = \sum_{k=1}^{K} W_k \sum_{m=1}^{M_2} A_{k,m} W_m Q_t(p + \Delta p_{km})$$

(3)

where $K$ represents the number of attention heads, $M_2$ indicates the number of sampling points, $W_k$ and $W_m$ are the learning weights, $A_{km}$ denotes the normalized attention weight, and $p + \Delta p_{km}$ signifies the learnable sample point position in 3D space. The feature is computed through trilinear interpolation from the voxel feature at this position.

**Iterative Structure Rectification Module.** While our initial octree structure is derived from information obtained through image segmentation, it may not precisely match the scene due to segmentation errors and projection inaccuracies. Simultaneously, the encoded octree query acquires valuable spatial information and the octree structure predicted from it complements and corrects the initial prediction, enabling us to address the shortcomings resulting from segmentation challenges.

Specifically, we partition the previous octree structure into two parts: the high-confidence region and the low-confidence region. By sorting the octree split probability values stored in the octree mask from largest to smallest and selecting the top K percent of the regions, we identify the location of high-confidence regions. These are areas where we believe the predictions are accurate and do not necessitate additional adjustments to prevent degradation.

For regions where confidence remains low, it is crucial to refine their outcomes by incorporating additional information sources. Given that predicting the octree structure requires a dense query format, we initially transform the octree query into a dense query, denoted as $Q_{dense} \in \mathbb{R}^{X \times Y \times Z \times C}$, using the procedure outlined in Section 3.3. Subsequently, we downsample it to each level of the octree through trilinear interpolation, resulting in $Q_{dense}^l \in \mathbb{R}^{\frac{X}{2^{L-l}} \times \frac{Y}{2^{L-l}} \times \frac{Z}{2^{L-l}} \times C}$. By utilizing the index of low-confidence regions, we can extract the corresponding features $Q_{lcr}^l \in \mathbb{R}^{N_l \times C}, l = 0, 1, \ldots, L - 1$. Employing a Multilayer Perceptron (MLP) to predict octree split probabilities from $Q_{lcr}^l$, we apply a weighted sum with the previous predictions of low-confidence regions to obtain refined predictions. These refined predictions are then concatenated with the preserved predictions of high-confidence areas, re-

sulting in the final rectified octree mask. The mask facilitates the transformation of the dense query into a rectified octree query for the subsequent layer.

### 3.6. Loss Function

To train the model, we use focal loss $L_{focal}$, lovasz-softmax loss $L_{ls}$, dice loss $L_{dice}$, affinity loss $L_{scal}^{geo}$ and $L_{scal}^{sem}$ from MonoScene[3]. In addition, we also use focal loss to supervise the octree prediction. The overall loss function $L = L_{focal} + L_{ls} + L_{dice} + L_{scal}^{geo} + L_{scal}^{sem} + L_{octree}$.

## 4. Experiments

In this section, we commence by introducing the datasets, evaluation metrics, and implementation details. Following this, we present the results and analysis from comprehensive comparison experiments and ablation studies to validate the efficacy and superiority of our method.

### 4.1. Datasets

**Occ3D-nuScenes[35]** constitutes a re-annotation of the nuScenes dataset[2], where precise ground truth labels for occupancy have been derived by aggregating LiDAR scans and human annotations. These annotations adhere to the ego coordinate system framework, consisting of 700 instances for training and 150 for validation. The spatial extents of occupancy along the X and Y axes in ego-coordinates are limited to the range of -40 meters to 40 meters, while the Z-axis spans from -1 meter to 5.4 meters. The occupancy labels are categorized into 17 distinct classes, with each class representing a volumetric space of 0.4 meters in each dimension (0.4m × 0.4m × 0.4m). Additionally, the dataset includes semantic labeling for the 18th category, denoted as "free", which represents the empty regions within a scene. Finally, Occ3D-nuScenes provides visibility masks for both LIDAR and camera modalities.

**SemanticKITTI[1]** comprises 22 distinct outdoor driving scenarios, with a focus on areas located in the forward trajectory of the vehicle. Each sample in this dataset covers a spatial extent ranging from [0.0m, -25.6m, -2.0m, 51.2m, 25.6m, 4.4m], with a voxel granularity set at [0.2m, 0.2m, 0.2m]. The dataset consists of volumetric representations, specifically in the form of 256×256×32 voxel grids. These grids undergo meticulous annotation with 21 distinct semantic classes. The voxel data is derived through a rigorous post-processing procedure applied to Lidar scans.

### 4.2. Evaluation metrics

Both 3D occupancy prediction and semantic scene completion utilize intersection-over-union (mIoU) over all classes as evaluation metrics, calculated as follows:

$$mIoU = \frac{1}{C} \sum_{c=1}^{C} \frac{TP_c}{TP_c + FP_c + FN_c},$$

(4)

Table 1. **3D Occupancy prediction performance** on the Occ3D-nuScenes *val* set. "⋆" denotes training with the camera mask.

| Method | Image Backbone | Published | mIoU | others | barrier | bicycle | bus | car | const. veh. | motorcycle | pedestrain | traffic cone | trailer | truck | drive. suf. | other flat | sidewalk | terrain | manmade | vegetation |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MonoScene[3] | EfficientNetB7 | CVPR 2022 | 6.06 | 1.75 | 7.23 | 4.26 | 4.93 | 9.38 | 5.67 | 3.98 | 3.01 | 5.90 | 4.45 | 7.17 | 14.91 | 6.32 | 7.92 | 7.43 | 1.01 | 7.65 |
| BEVDet[8] | ResNet101 | arxiv 2021 | 11.73 | 2.09 | 15.29 | 0.0 | 4.18 | 12.97 | 1.35 | 0.0 | 0.43 | 0.13 | 6.59 | 6.66 | 52.72 | 19.04 | 26.45 | 21.78 | 14.51 | 15.26 |
| BEVFormer[18] | ResNet101 | ECCV 2022 | 23.67 | 5.03 | 38.79 | 9.98 | 34.41 | 41.09 | 13.24 | 16.50 | 18.15 | 17.83 | 18.66 | 27.70 | 48.95 | 27.73 | 29.08 | 25.38 | 15.41 | 14.46 |
| BEVStereo[15] | ResNet101 | AAAI 2023 | 24.51 | 5.73 | 38.41 | 7.88 | 38.70 | 41.20 | 17.56 | 17.33 | 14.69 | 10.31 | 16.84 | 29.62 | 54.08 | 28.92 | 32.68 | 26.54 | 18.74 | 17.49 |
| TPVFormer[9] | ResNet101 | CVPR 2023 | 28.34 | 6.67 | 39.20 | 14.24 | 41.54 | 46.98 | 19.21 | 22.64 | 17.87 | 14.54 | 30.20 | 35.51 | 56.18 | 33.65 | 35.69 | 31.61 | 19.97 | 16.12 |
| OccFormer[45] | ResNet101 | ICCV 2023 | 21.93 | 5.94 | 30.29 | 12.32 | 34.40 | 39.17 | 14.44 | 16.45 | 17.22 | 9.27 | 13.90 | 26.36 | 50.99 | 30.96 | 34.66 | 22.73 | 6.76 | 6.97 |
| CTF-Occ[35] | ResNet101 | arxiv 2023 | 28.53 | 8.09 | 39.33 | 20.56 | 38.29 | 42.24 | 16.93 | 24.52 | 22.72 | 21.05 | 22.98 | 31.11 | 53.33 | 33.84 | 37.98 | 33.23 | 20.79 | 18.00 |
| RenderOcc[27] | ResNet101 | arxiv 2023 | 26.11 | 4.84 | 31.72 | 10.72 | 27.67 | 26.45 | 13.87 | 18.2 | 17.67 | 17.84 | 21.19 | 23.25 | 63.20 | 36.42 | 46.21 | 44.26 | 19.58 | 20.72 |
| BEVDet4D[7]* | Swin-B | arxiv 2022 | 42.02 | **12.15** | 49.63 | 25.1 | 52.02 | 54.46 | 27.87 | 27.99 | 28.94 | 27.23 | 36.43 | 42.22 | 82.31 | 43.29 | 54.46 | 57.9 | 48.61 | 43.55 |
| PanoOcc[41]* | ResNet101 | arxiv 2023 | 42.13 | 11.67 | 50.48 | 29.64 | 49.44 | 55.52 | 23.29 | **33.26** | 30.55 | 30.99 | 34.43 | 42.57 | 83.31 | 44.23 | 54.40 | 56.04 | 45.94 | 40.40 |
| FB-OCC[19]* | ResNet101 | ICCV 2023 | 43.41 | 12.10 | 50.23 | **32.31** | 48.55 | 52.89 | **31.20** | 31.25 | **30.78** | **32.33** | 37.06 | 40.22 | **83.34** | **49.27** | **57.13** | **59.88** | 47.67 | 41.76 |
| Ours* | ResNet101 | - | **44.02** | 11.96 | **51.70** | 29.93 | **53.52** | **56.77** | 30.83 | 33.17 | 30.65 | 29.99 | **37.76** | **43.87** | 83.17 | 44.52 | 55.45 | 58.86 | **49.52** | 46.33 |

Table 2. **Comparison of query form and efficiency** on the Occ3D-nuScenes *val* set. The experiments were conducted on the NVIDIA A100 GPU.

| Methods | Query Form | Query Resolution | mIoU | Latency | Memory |
|---|---|---|---|---|---|
| BEVFormer[18] | 2D BEV | 200×200 | 23.67 | 302 ms | 25100M |
| TPVFormer[9] | 2D tri-plane | 200×(200+16+16) | 28.34 | 341 ms | 29000M |
| PanoOcc[41] | 3D voxel | 100×100 ×16 | 42.13 | 502 ms | 35000M |
| FBOCC[19] | 3D voxel & 2D BEV | 200×200×16 & 200×200 | 43.41 | 463 ms | 31000M |
| Ours | Octree Query | - | **44.02** | 386 ms | 26500M |

where $TP_c$, $FP_c$, and $FN_c$ correspond to the number of true positive, false positive, and false negative predictions for class $c_i$, and $C$ is the number of classes.

### 4.3. Implementation Details

Based on previous research, we set the input image size to 900×1600 and employ ResNet101-DCN[5] as the image backbone. We extract multi-scale features from the Feature Pyramid Network (FPN)[20] with downsample sizes of 1/8, 1/16, 1/32, and 1/64. The feature dimension $C$ is set to 256. The octree depth is 3, and the initial query resolution is 50×50×4. For $\alpha$, we designate $\alpha^1$ as 0.2 and $\alpha^2$ as 0.6. The octree encoder comprises three layers, each composed of Temporal Self-Attention (TSA), Image Cross-Attention (ICA), and Iterative Structure Rectification (ISR) modules. Both $M_1$ and $M_2$ are set to 4. In TSA, we fuse 4 frames with a time interval of 0.5s. In ISR, the top 10% predictions are considered high-confidence in level 1, and 30% in level 2. The loss weights are uniformly set to 1.0. For optimization, we employ Adam[12] optimizer with a learning rate of 2e-4 and weight decay of 0.01. The batch size is 8, and the model is trained for 24 epochs, consuming 3 days on 8 NVIDIA A100 GPUs.

### 4.4. Main Results

**3D Occupancy Prediction** In Tab. 1, we provide comparasion experiments with other state-of-the-art occupancy prediction methodds on Occ3d-nus validation set. The performance of FBOCC[19] relies on open-source code, which we evaluate after ensuring consistency in details (utilizing the same backbone, image resolution, and excluding CBGS[47]) for a fair comparison. Performance for other methods are reported in a series of works[27, 35, 41]. Our approach demonstrates superior performance on mIoU compared to previous methods, particularly excelling in foreground classes such as barriers, cars, and buses, as well as in scene structure classes like manmade and vegetation. This highlights that processing scenes using multiple granularities aligns better with the scene characteristics, enhancing overall expressiveness.

Moreover, we evaluate the efficacy of our approach by comparing it to alternative methods utilizing diverse query forms, including 2D BEV query, 2D tri-plane query, and 3D query, as depicted in Tab 2. The results indicate that our approach not only surpasses these methods in terms of performance but also demonstrates significantly reduced memory usage and lower latency compared to dense queries, approaching the levels observed with 2D queries. Clearly, sparse octree queries effectively mitigate computational overhead while ensuring robust performance.

**3D Semantic Scene Completion.** To better assess the effectiveness of our approach, we conduct comparative experiments for the Semantic Scene Completion (SSC) task. As demonstrated in Tab 3, we compare our results with those of other SSC methods on the SemanticKITTI validation set. The outcomes indicate that, owing to the construction and correction of the octree, our model exhibits a more accurate perception of space. Our method outperforms other methods in the IoU metric, which evaluates the quality of reconstructed geometry. Additionally, for specific semantic classes such as car and vegetation, we achieve superior results, attributed to the enhanced treatment of objects facilitated by multi-granularity modeling.

### 4.5. Ablation Study & Analysis

In this subsection, we perform ablation studies on the Occ3d-nus validation set to assess the effectiveness of our proposed modules. Additionally, we analyze the impact of octree depth and the number of queries at each octree level.

Table 3. **3D Semantic Scene Completion** performance on SemanticKITTI *val* set.

| Method | Published | IoU | mIoU | road | sidewalk | parking | other-ground | building | car | truck | bicycle | motorcycle | other-vehicle | vegetation | trunk | terrain | person | bicyclist | motorcyclist | fence | pole | traf.-sign |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| LMSCNet[31] | 3DV 2020 | 28.61 | 6.70 | 40.68 | 18.22 | 4.38 | 0.00 | 10.31 | 18.33 | 0.00 | 0.00 | 0.00 | 0.00 | 13.66 | 0.02 | 20.54 | 0.00 | 0.00 | 0.00 | 1.21 | 0.00 | 0.00 |
| AICNet[14] | CVPR 2020 | 29.59 | 8.31 | 43.55 | 20.55 | 11.97 | 0.07 | 12.94 | 14.71 | 4.53 | 0.00 | 0.00 | 0.00 | 15.37 | 2.90 | 28.71 | 0.00 | 0.00 | 0.00 | 2.52 | 0.06 | 0.00 |
| 3DSketch[4] | CVPR 2020 | 33.30 | 7.50 | 41.32 | 21.63 | 0.00 | 0.00 | 14.81 | 18.59 | 0.00 | 0.00 | 0.00 | 0.00 | 19.09 | 0.00 | 26.40 | 0.00 | 0.00 | 0.00 | 0.73 | 0.00 | 0.00 |
| JS3C-Net[43] | AAAI 2021 | 38.98 | 10.31 | 50.49 | 23.74 | 11.94 | 0.07 | 15.03 | 24.65 | 4.41 | 0.00 | 0.00 | 6.15 | 18.11 | 4.33 | 26.86 | 0.67 | 0.27 | 0.00 | 3.94 | 3.77 | 1.45 |
| MonoScene[3] | CVPR 2022 | 36.86 | 11.08 | 56.52 | 26.72 | 14.27 | 0.46 | 14.09 | 23.26 | 6.98 | 0.61 | 0.45 | 1.48 | 17.89 | 2.81 | 29.64 | 1.86 | 1.20 | 0.00 | 5.84 | 4.14 | 2.25 |
| TPVFormer[9] | CVPR 2023 | 35.61 | 11.36 | 56.50 | 25.87 | 20.60 | 0.85 | 13.88 | 23.81 | 8.08 | 0.36 | 0.05 | 4.35 | 16.92 | 2.26 | 30.38 | 0.51 | 0.89 | 0.00 | 5.94 | 3.14 | 1.52 |
| VoxFormer[17] | CVPR 2023 | 44.02 | 12.35 | 54.76 | 26.35 | 15.50 | 0.70 | 17.65 | 25.79 | 5.63 | 0.59 | 0.51 | 3.77 | 24.39 | 5.08 | 29.96 | 1.78 | 3.32 | 0.00 | 7.64 | 7.11 | 4.18 |
| OccFormer[45] | ICCV 2023 | 36.50 | 13.46 | 58.84 | 26.88 | 19.61 | 0.31 | 14.40 | 25.09 | 25.53 | 0.81 | 1.19 | 8.52 | 19.63 | 3.93 | 32.63 | 2.78 | 2.82 | 0.00 | 5.61 | 4.26 | 2.86 |
| Symphonies[10] | arxiv 2023 | 41.44 | 13.44 | 55.78 | 26.77 | 14.57 | 0.19 | 18.76 | 27.23 | 15.99 | 1.44 | 2.28 | 9.52 | 24.50 | 4.32 | 28.49 | 3.19 | 8.09 | 0.00 | 6.18 | 8.99 | 5.39 |
| Ours | - | 44.71 | 13.12 | 55.13 | 26.74 | 18.68 | 0.65 | 18.69 | 28.07 | 16.43 | 0.64 | 0.71 | 6.03 | 25.26 | 4.89 | 32.47 | 2.25 | 2.57 | 0.00 | 4.01 | 3.72 | 2.36 |

Table 4. **Ablation experiments of Modules** on Occ3d-nuScenes dataset with reducing the input image size to 0.3x. "Sem.Init." demotes semantic-guided octree initialization. "Iter.Rec." means iterative structure rectification. The experiments were conducted on the NVIDIA A40 GPU.

| Baseline | Octree Query | Sem.Init. | Iter.Rec. | mIoU | Latency | Memory |
|---|---|---|---|---|---|---|
| ✓ | | | | 34.17 | 266 ms | 27200M |
| ✓ | ✓ | | | 34.91 | 218 ms | 18300M |
| ✓ | ✓ | ✓ | | 36.63 | 214 ms | 17900M |
| ✓ | ✓ | | ✓ | 35.88 | 227 ms | 18500M |
| ✓ | ✓ | ✓ | ✓ | **37.40** | 224 ms | 18500M |

Table 5. **Ablation for different octree depth** on Occ3d-nuScenes *val* set with reducing the input image size to 0.3x. The experiments were conducted on the NVIDIA A40 GPU.

| | Query Form | Query Resolution | mIoU | Latency | Memory |
|---|---|---|---|---|---|
| (a) | | 50×50×4 | 28.51 | 129 ms | 7400M |
| (b) | 3D Voxel | 50×50×16 | 31.67 | 186 ms | 11600M |
| (c) | | 100×100×8 | 32.21 | 204 ms | 17400M |
| (d) | | 100×100×16 | 34.10 | 266 ms | 27200M |
| (e) | | 50×50×4 to 100×100×8 | 32.02 | 182 ms | 12400M |
| (f) | Octree Query | 50×50×4 to 100×100×16 | 33.76 | 207 ms | 16800M |
| (g) | | 50×50×4 to 200×200×16 | **37.40** | 224 ms | 18500M |

Table 6. **Ablation for the choice of query selection ratio** on Occ3d-nuScenes *val* set with reducing the input image size to 0.3x. The experiments were conducted on the NVIDIA A40 GPU.

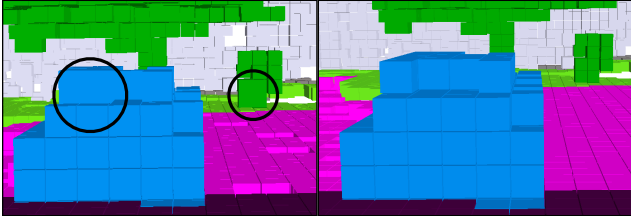| | Query Form | Query Selection Ratio | mIoU | Latency | Memory |
|---|---|---|---|---|---|
| (a)) | Octree Query | 25%, 50% | 36.73 | 220 ms | 18000M |
| (b) | Octree Query | 25%, 60% | 36.12 | 255 ms | 21000M |
| (c) | Octree Query | 20%, 60% | **37.10** | 224 ms | 18500M |



Figure 3. **Illustration of octree structure rectification.** The left figure displays the initially predicted octree structure, while the right figure depicts the octree structure after Iterative Structure Rectification module. It is evident that the predicted octree structure becomes more consistent with the object's shape following the rectification module.

**Effectiveness of Octree Queries.** To validate the superiority of octree queries, we maintain consistent TSA and ICA settings while removing the proposed octree structure initialization and rectification modules. This allows for a comparison with the baseline that employs dense queries of size 100×100×16 (the largest size achievable within memory limitations). As illustrated in Table 4, the experimental results reveal that we continue to outperform with a 0.8 mIoU advantage, even with a memory saving of approximately 9G. This demonstrates that octree prediction can adeptly assign queries with different granularities to distinct semantic regions. The constructed octree queries are accommodate to various object shapes, showcasing a more efficient utilization of computational budget.

**Effectiveness of Semantic-guided Structure Initialization module.** To underscore the significance of the initial octree structure, we exclude the Semantic-guided Octree Initialization module and employ a MLP to predict the octree structure from randomly initialized queries. This results in a performance decrease of 1.7 mIoU, highlighting that despite having ground truth for supervision, the structurally predicted information from the initialized query is inaccurate due to the absence of coding for obtaining valid information. The incorporation of semantic priors proves crucial as it serves as a valuable reference to enhance the quality of the initialized octree, thereby improving overall model performance.

**Effectiveness of Iteritive Structure Rectification module.** We additionally perform an ablation study on the Iterative Structure Rectification module, as shown in Table 4. The incorporation of this module has led to noticeable improvements in performance. Additionally, we have presented a visualization of the octree structure both before and after correction (Fig. 3). These visualizations highlight the

| CAM_FRONT_LEFT | CAM_FRONT | CAM_FRONT_RIGHT | CAM_BACK_LEFT | CAM_BACK | CAM_BACK_RIGHT |

PanoOcc          FBOCC          Ours          GT

Barrier  Bicycle  Bus  Car  Const. Veh.  Motor  Ped.  Traf. C.  Trailer  Truck  Drive. Surf.  Flat  Sidewalk  Terrain  Manmade  Vegetation
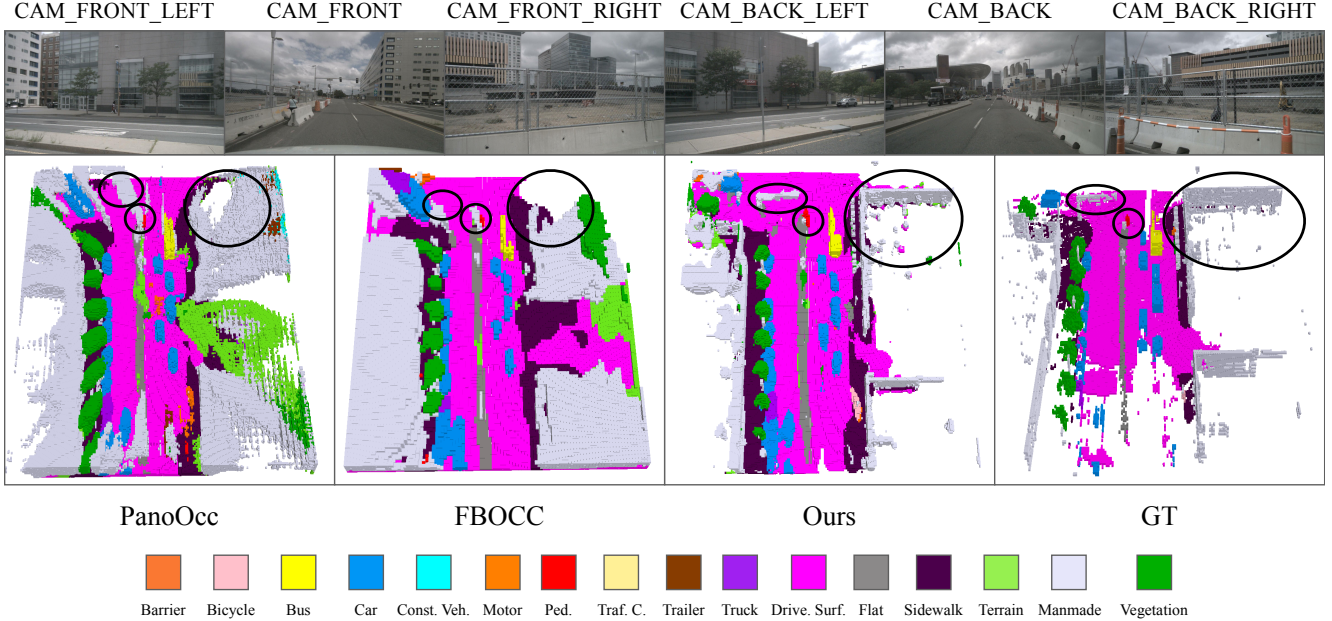
Figure 4. **Qualitative results on Occ3D-nuScenes validation set.** The first row displays input multi-view images, while the second row showcases the occupancy prediction results of PanoOcc[41], FBOCC[19], our methods, and the ground truth.

module's role in rectifying inaccuracies in the initialized octree structure, thereby enhancing its compatibility with input scenarios. Consequently, this refinement contributes to the efficiency of octree query expression, positively impacting overall performance.

**Discussion on the depth of octree.** As shown in Tab. 5, we conducted experiments to investigate the impact of variations in the depth of the constructed octree. (a), (b), (c), and (d) illustrate the performance of four different sizes of 3D queries, while (e), (f), and (g) depict the performance of octree queries with three different depths and initial query sizes. By comparing (e) with (c) and (f) with (d), it is evident that our octree query approach achieves performance comparable to dense queries while significantly reducing memory consumption. (g) demonstrates that we can further refine modeling granularity, surpassing the constraints imposed by the excessive graphics memory overhead (>80G) associated with the $200\times200\times16$ dense modeling and acheving higher performance.

**Discussion on query selection ratio of each level in the octree.** In Table 6, we present the results obtained by maintaining various query ratios across different levels of the octree. A comparison between (a) and (b) reveals that an excessive number of fine-grained queries leads to a degradation in performance. This is primarily due to the fact that the scenario does not predominantly consist of foreground points requiring modeling at such a detailed level. Further comparison between (a) and (c) demonstrates that, given the prevalence of coarsest-grained semantic regions, par-

ticularly those containing mostly empty spaces, modeling this portion at a fine-grained level adversely impacts computational efficiency. Consequently, determining the optimal number of queries with different granularities should align with the statistical distribution of object granularity within the scene. This ensures that queries of varying granularity can effectively handle semantic regions corresponding to their specific granularity levels.

### 4.6. Visualization

The Fig. 4 displays qualitative results obtained by our methods and other occupancy prediction solutions on the Occ3D-nus validation set, where the resolution of the voxel predictions is $200\times200\times16$. The figure illustrates that our approach comprehensively understands the structure of the scene, showcasing superior performance in scene understanding.

### 5. Conclusions

In conclusion, our paper introduces OctreeOcc, a novel 3D occupancy prediction framework that addresses the limitations of traditional dense-grid representations in understanding 3D scenes. OctreeOcc's adaptive utilization of octree representations enables the capture of valuable information with variable granularity, catering to objects of diverse sizes and complexities. Our extensive experimental results affirm OctreeOcc's capability to attain state-of-the-art performance in 3D occupancy prediction while concurrently reducing computational overhead.

# References

[1] 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 *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9297–9307, 2019. 5

[2] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020. 5

[3] Anh-Quan Cao and Raoul de Charette. Monoscene: Monocular 3d semantic scene completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3991–4001, 2022. 2, 5, 6, 7

[4] 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. 7

[5] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 764–773, 2017. 6

[6] Chunyang Fu, Ge Li, Rui Song, Wei Gao, and Shan Liu. Octattention: Octree-based large-scale contexts model for point cloud compression. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 625–633, 2022. 3

[7] Junjie Huang and Guan Huang. Bevdet4d: Exploit temporal cues in multi-camera 3d object detection. *arXiv preprint arXiv:2203.17054*, 2022. 6

[8] 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, 6

[9] Yuanhui Huang, Wenzhao Zheng, Yunpeng Zhang, Jie Zhou, and Jiwen Lu. Tri-perspective view for vision-based 3d semantic occupancy prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9223–9232, 2023. 1, 2, 6, 7

[10] Haoyi Jiang, Tianheng Cheng, Naiyu Gao, Haoyang Zhang, Wenyu Liu, and Xinggang Wang. Symphonize 3d semantic scene completion with contextual instance queries. *arXiv preprint arXiv:2306.15670*, 2023. 1, 2, 7

[11] Yanqin Jiang, Li Zhang, Zhenwei Miao, Xiatian Zhu, Jin Gao, Weiming Hu, and Yu-Gang Jiang. Polarformer: Multicamera 3d object detection with polar transformer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1042–1050, 2023. 2

[12] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. 6

[13] Huan Lei, Naveed Akhtar, and Ajmal Mian. Octree guided cnn with spherical kernels for 3d point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9631–9640, 2019. 3

[14] 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. 7

[15] Yinhao Li, Han Bao, Zheng Ge, Jinrong Yang, Jianjian Sun, and Zeming Li. Bevstereo: Enhancing depth estimation in multi-view 3d object detection with dynamic temporal stereo, 2022. 2, 6

[16] 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. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1477–1485, 2023. 2

[17] Yiming Li, Zhiding Yu, Christopher Choy, Chaowei Xiao, Jose M. Alvarez, Sanja Fidler, Chen Feng, and Anima Anandkumar. Voxformer: Sparse voxel transformer for camera-based 3d semantic scene completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9087–9098, 2023. 1, 2, 7

[18] Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng Dai. Bevformer: Learning bird's-eye-view representation from multi-camera images via spatiotemporal transformers. In *European conference on computer vision*, pages 1–18. Springer, 2022. 2, 6

[19] Zhiqi Li, Zhiding Yu, Wenhai Wang, Anima Anandkumar, Tong Lu, and Jose M Alvarez. Fb-bev: Bev representation from forward-backward view transformations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6919–6928, 2023. 1, 6, 8, 2

[20] 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. 6

[21] Xuewu Lin, Tianwei Lin, Zixiang Pei, Lichao Huang, and Zhizhong Su. Sparse4d: Multi-view 3d object detection with sparse spatial-temporal fusion. *arXiv preprint arXiv:2211.10581*, 2022. 2

[22] Yingfei Liu, Tiancai Wang, Xiangyu Zhang, and Jian Sun. Petr: Position embedding transformation for multi-view 3d object detection. In *European Conference on Computer Vision*, pages 531–548. Springer, 2022.

[23] Yingfei Liu, Junjie Yan, Fan Jia, Shuailin Li, Aqi Gao, Tiancai Wang, and Xiangyu Zhang. Petrv2: A unified framework for 3d perception from multi-camera images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3262–3272, 2023. 2

[24] Zhijian Liu, Haotian Tang, Alexander Amini, Xinyu Yang, Huizi Mao, Daniela L Rus, and Song Han. Bevfusion: Multi-task multi-sensor fusion with unified bird's-eye view representation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2774–2781. IEEE, 2023. 2

[25] Donald Meagher. Geometric modeling using octree encod-

ing. *Computer graphics and image processing*, 19(2):129–147, 1982. 2, 3

[26] Ruihang Miao, Weizhou Liu, Mingrui Chen, Zheng Gong, Weixin Xu, Chen Hu, and Shuchang Zhou. Occdepth: A depth-aware method for 3d semantic scene completion. *arXiv preprint arXiv:2302.13540*, 2023. 2

[27] Mingjie Pan, Jiaming Liu, Renrui Zhang, Peixiang Huang, Xiaoqi Li, Li Liu, and Shanghang Zhang. Renderocc: Vision-centric 3d occupancy prediction with 2d rendering supervision, 2023. 2, 6

[28] Jinhyung Park, Chenfeng Xu, Shijia Yang, Kurt Keutzer, Kris Kitani, Masayoshi Tomizuka, and Wei Zhan. Time will tell: New outlooks and a baseline for temporal multi-view 3d object detection. *arXiv preprint arXiv:2210.02443*, 2022. 2

[29] Jonah Philion and Sanja Fidler. Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16*, pages 194–210. Springer, 2020. 2

[30] Zizheng Que, Guo Lu, and Dong Xu. Voxelcontext-net: An octree based framework for point cloud compression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6042–6051, 2021. 3

[31] 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

[32] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pages 234–241. Springer, 2015. 4

[33] Chonghao Sima, Wenwen Tong, Tai Wang, Li Chen, Silei Wu, Hanming Deng, Yi Gu, Lewei Lu, Ping Luo, Dahua Lin, and Hongyang Li. Scene as occupancy. 2023. 2, 4

[34] Zhiyu Tan, Zichao Dong, Cheng Zhang, Weikun Zhang, Hang Ji, and Hao Li. Ovo: Open-vocabulary occupancy. *arXiv preprint arXiv:2305.16133*, 2023. 1, 2

[35] 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. 2, 5, 6

[36] Peng-Shuai Wang. Octformer: Octree-based transformers for 3d point clouds. *arXiv preprint arXiv:2305.03045*, 2023. 3

[37] Peng-Shuai Wang, Yang Liu, Yu-Xiao Guo, Chun-Yu Sun, and Xin Tong. O-cnn: Octree-based convolutional neural networks for 3d shape analysis. *ACM Transactions On Graphics (TOG)*, 36(4):1–11, 2017. 3

[38] Peng-Shuai Wang, Chun-Yu Sun, Yang Liu, and Xin Tong. Adaptive o-cnn: A patch-based deep representation of 3d shapes. *ACM Transactions on Graphics (TOG)*, 37(6):1–11, 2018. 3

[39] Shihao Wang, Yingfei Liu, Tiancai Wang, Ying Li, and Xiangyu Zhang. Exploring object-centric temporal modeling for efficient multi-view 3d object detection. *arXiv preprint arXiv:2303.11926*, 2023. 2

[40] 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. 2

[41] Yuqi Wang, Yuntao Chen, Xingyu Liao, Lue Fan, and Zhaoxiang Zhang. Panoocc: Unified occupancy representation for camera-based 3d panoptic segmentation. *arXiv preprint arXiv:2306.10013*, 2023. 1, 2, 4, 6, 8

[42] Yi Wei, Linqing Zhao, Wenzhao Zheng, Zheng Zhu, Jie Zhou, and Jiwen Lu. Surroundocc: Multi-camera 3d occupancy prediction for autonomous driving. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 21729–21740, 2023. 1, 2

[43] 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*, pages 3101–3109, 2021. 7

[44] Chenyu Yang, Yuntao Chen, Hao Tian, Chenxin Tao, Xizhou Zhu, Zhaoxiang Zhang, Gao Huang, Hongyang Li, Yu Qiao, Lewei Lu, et al. Bevformer v2: Adapting modern image backbones to bird's-eye-view recognition via perspective supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17830–17839, 2023. 2

[45] 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

[46] Chao Zhou, Yanan Zhang, Jiaxin Chen, and Di Huang. Octr: Octree-based transformer for 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5166–5175, 2023. 3

[47] 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

[48] 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. 4

# OctreeOcc: Efficient and Multi-Granularity Occupancy Prediction Using Octree Queries

## Supplementary Material

## A. More Details

In this section, we provide detailed explanations of our proposed modules.

For **Semantic-Guided Octree Initialization**, we started by obtaining semantic segmentation labels for images through the projection of occupancy labels onto the surround-view image. Subsequently, a UNet semantic segmentation model was trained using these labels. To initialize the octree structure and queries, we conducted offline generation of semantic segmentation result graphs. This offline approach helps avoid unfair comparisons that may arise from simultaneous training with the occupancy prediction model. The initialization process involves the random initialization of dense queries and the coordinates of each query's center point serve as reference points, which are projected onto the range-view image. If projected onto a ground pixel (representing driveable surface, other flat, or sidewalk), the probability increases by 0.1. Conversely, if projected onto a background pixel (excluding ground classes), the probability increases by 0.5. If a reference point is projected onto a foreground pixel, the probability of requiring a split at that position increases by 1.0. This process assigns a split probability to each query, and the octree mask is then constructed through average pooling, capturing split probabilities at different query levels. After obtaining the octree mask, we designate the top 20% confidence queries as parent queries in level 1, while the remaining queries become Leaf queries and remain unsplit. Moving to level 2, after splitting the parent queries into octants, the top 60% confidence positions are selected as new parent queries, and the remainder as leaf queries. By storing leaf queries at each level, we construct a sparse and multi-granularity octree structure for queries.

For **Iterative Structure Rectification**, in level 1, we retain predictions for the top 10% of positions with confidence. For the remaining positions, we use a 2-layer MLP to predict probabilities. The new probabilities are then calculated by blending them with the existing probabilities, with a weight distribution of 60% for the new probabilities and 40% for the old ones. We identify the top 10% of positions with the new probability values as the required splits, constituting the structure of the new level 1. Similarly, in level 2, we preserve predictions for the top 30% of positions with confidence. For the positions not in the top 30%, we predict probabilities using a 2-layer MLP. The new probabilities are computed by combining them with the original probabilities, with a weight distribution of 50% for

each. The top 30% of the new probability values are then selected as the positions that need to be split, forming the structure of the new level 2.

Table 7. Comparison of query form and efficiency on the SemanticKITTI *val* set. The experiments were conducted on the NVIDIA A40 GPU.

| Methods | IoU | mIoU | Latency | Memory |
|---|---|---|---|---|
| TPVFormer[9] | 35,61 | 11,36 | 179 ms | 23000M |
| VoxFormer[17] | 44.02 | 12.35 | 177 ms | 23200M |
| OccFormer[45] | 36.50 | 13.46 | 173 ms | 22400M |
| Symphonics[10] | 41.44 | 13.44 | 187 ms | 22000M |
| Ours | 44.71 | 13.12 | **162 ms** | **19000M** |

Table 8. Comparison of octree quality at different stages

| Stage | mIoU(L1) | mIoU(L2) |
|---|---|---|
| Initialized | 57.34 | 51.28 |
| After the 1st Rectification | 60.13 | 53.95 |
| After the 2nd Rectification | 62.27 | 56.79 |

## B. Comparison of efficiency on the SemanticKITTI dataset

As shown in Tab. 7, we evaluate the efficacy of our approach by comparing it to alternative methods on the SemanticKITTI dataset. The results demonstrate that our approach not only outperforms these methods in terms of performance but also requires significantly less memory and exhibits lower latency in comparison. It is evident that employing sparse octree queries effectively mitigates computational overhead while maintaining robust performance.

## C. Discussion on Octree Structure Quality

In Tab. 8, we present an analysis of the influence of the proposed initialization and correction modules on the octree's quality. It is observed that continuous correction of the octree structure throughout the network learning process resulted in the rectification of several inaccurate predictions (approximately a 10% increase in mIoU predicted at each level).

CAM_FRONT_LEFT  CAM_FRONT  CAM_FRONT_RIGHT  CAM_BACK_LEFT  CAM_BACK  CAM_BACK_RIGHT

PanoOcc  FBOCC  Ours  GT

Barrier  Bicycle  Bus  Car  Const. Veh.  Motor  Ped.  Traf. C.  Trailer  Truck  Drive. Surf.  Flat  Sidewalk  Terrain  Manmade  Vegetation
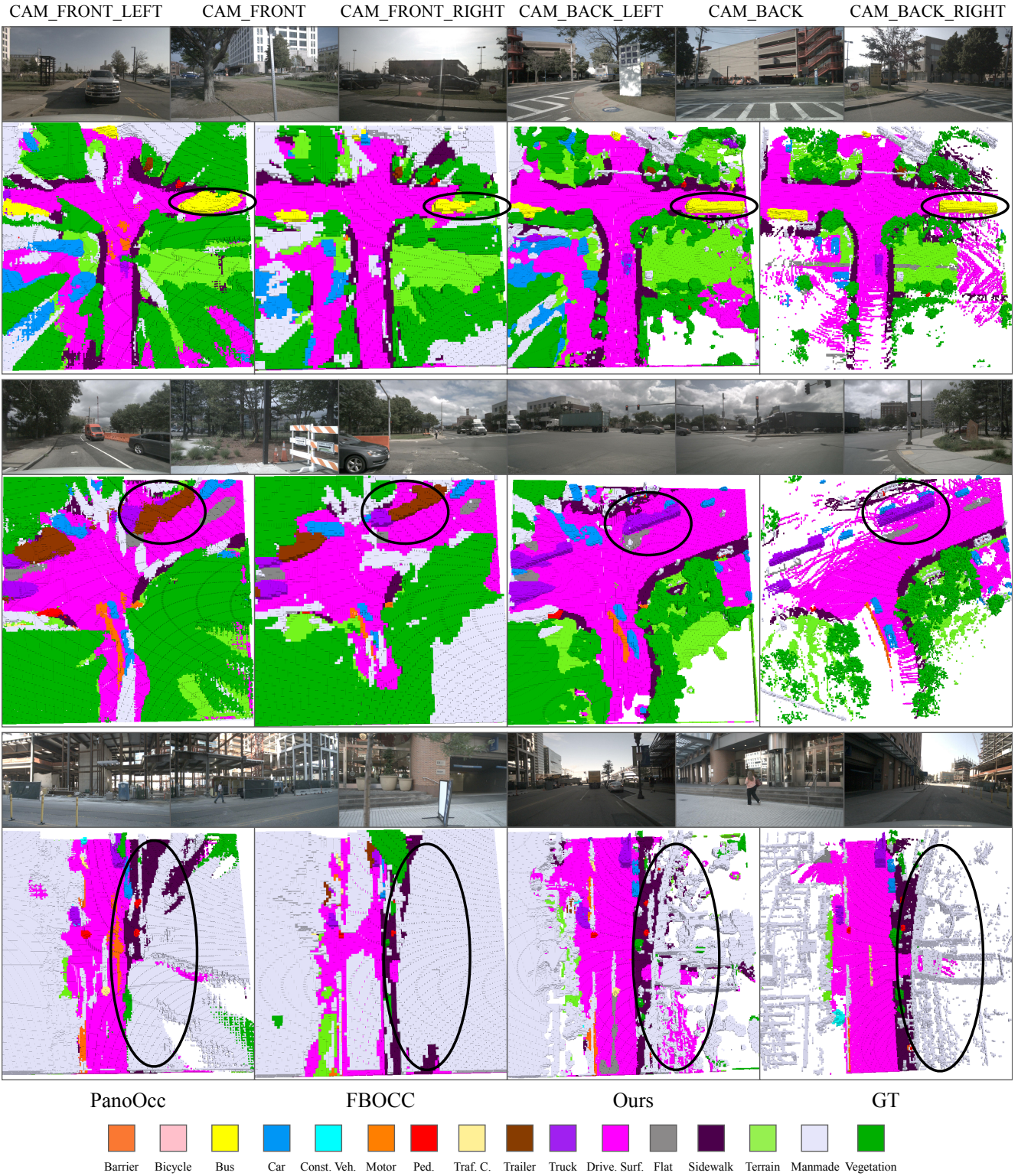
Figure 5. **More visualization on Occ3D-nuScenes validation set.** The first row displays input multi-view images, while the second row showcases the occupancy prediction results of PanoOcc[41], FBOCC[19], our methods, and the ground truth.

## D. More Visualization

Fig. 5 shows additional visualizations of proposed Oc-treeOcc. Evidently, our approach, leveraging the multi-granularity octree modeling, demonstrates superior perfor-mance particularly in the categories of truck, bus, and man-made objects.