# OPUS: Occupancy Prediction Using a Sparse Set

**Jiabao Wang**[1,*] **Zhaojiang Liu**[2*]**, Qiang Meng**[3]**, Liujiang Yan**[3]**, Ke Wang**[3]**, Jie Yang**[3]**,
**Wei Liu**[2]**, Qibin Hou**[1,4†]**, Ming-Ming Cheng**[1,4]
[1]VCIP, College of Computer Science, Nankai University
[2]Shanghai Jiao Tong University    [3]KargoBot Inc.    [4]NKIARI, Shenzhen Futian
https://github.com/jbwang1997/OPUS

## Abstract

Occupancy prediction, aiming at predicting the occupancy status within voxelized 3D environment, is quickly gaining momentum within the autonomous driving community. Mainstream occupancy prediction works first discretize the 3D environment into voxels, then perform classification on such dense grids. However, inspection on sample data reveals that the vast majority of voxels is unoccupied. Performing classification on these empty voxels demands suboptimal computation resource allocation, and reducing such empty voxels necessitates complex algorithm designs. To this end, we present a novel perspective on the occupancy prediction task: formulating it as a streamlined set prediction paradigm without the need for explicit space modeling or complex sparsification procedures. Our proposed framework, called OPUS, utilizes a transformer encoder-decoder architecture to simultaneously predict occupied locations and classes using a set of learnable queries. Firstly, we employ the Chamfer distance loss to scale the set-to-set comparison problem to unprecedented magnitudes, making training such model end-to-end a reality. Subsequently, semantic classes are adaptively assigned using nearest neighbor search based on the learned locations. In addition, OPUS incorporates a suite of non-trivial strategies to enhance model performance, including coarse-to-fine learning, consistent point sampling, and adaptive re-weighting, *etc*. Finally, compared with current state-of-the-art methods, our lightest model achieves superior RayIoU on the Occ3D-nuScenes dataset at near $2\times$ FPS, while our heaviest model surpasses previous best results by 6.1 RayIoU.

## 1 Introduction

Compared with well-established box representations [8, 23, 20, 39, 29, 48, 32], voxel based occupancy [16, 37, 10, 38, 33] can provide finer geometry and semantic information for the surrounding scene. For example, it is not straightforward to use bounding boxes to describe vehicles with doors open or cranes with outriggers deployed. While occupancy can naturally describe such uncommon shapes. Thus occupancy prediction is quickly gaining traction in the autonomous driving community.

Recent approaches [3, 46, 9, 27, 16, 34] to the task predominantly rely on dense data representation, with a direct one-to-one correspondence between feature points and physical voxels. It has come to our attention that the vast majority of physical voxels is empty. For instance, in SemanticKITTI [1], approximately 67% of all voxels are empty, while in Occ3D-nuScenes [38], this proportion exceeds 90%. Such sparse nature of occupancy data renders the direct dense representation undeniably inefficient, as majority of the computation is allocated towards empty voxels. Alternative sparse latent representations have been explored to alleviate such inefficiency, such as the Tri-Perspective View representation [37, 9] or reduced solution spaces [21, 10], leading to notably reduced computational costs. However, these approaches still treat occupancy prediction as a classification problem at specific locations, necessitating complex intermediate designs and explicit modeling of 3D spaces.

---

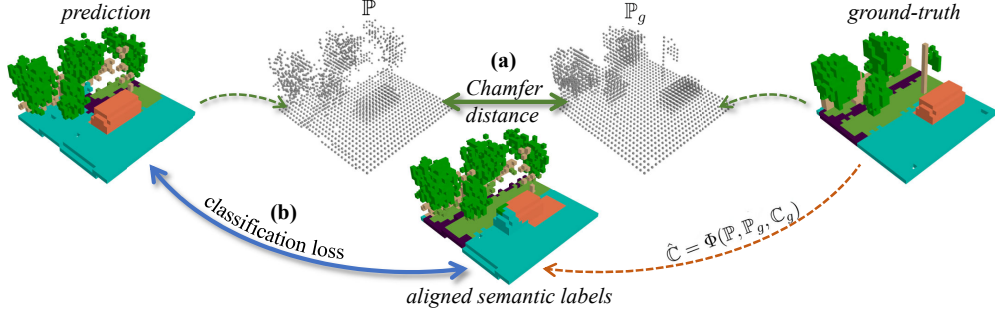[*]Equal contribution. [†]Corresponding author.

Figure 1: The occupancy prediction is approached as a set prediction problem. For each scene, we predict a set of point positions $\mathbb{P}$ and a set of the corresponding semantic classes $\mathbb{C}$. With the ground-truth set of occupied voxel positions $\mathbb{P}_g$ and classes $\mathbb{C}_g$, we decouple the set-to-set matching task into two distinct components: (a) Enforcing similarity in the point distributions of $\mathbb{P}$ and $\mathbb{P}_g$ using the Chamfer distance. (b) Aligning the predicted classes $\mathbb{C}$ with the ground-truths $\hat{\mathbb{C}} = \Phi(\mathbb{P}, \mathbb{P}_g, \mathbb{C}_g)$, where $\Phi$ generates a set of classes for points $\mathbb{P}$ based on those of the nearest ground-truth points.

In this work, we instead formulate the task as a direct set prediction problem, where we regress occupied locations and classify corresponding semantic labels in parallel. Our proposed framework termed OPUS leverages a transformer encoder-decoder architecture featuring: (1) an image encoder to extract 2D features from multi-view images; (2) a set of learnable queries to predict occupied locations and semantic classes; (3) a sparse decoder to update query features with correlated image features. Our OPUS eliminates the need for explicit space modeling or complex sparsification procedures, offering a streamlined and elegant end-to-end solution. However, a key challenge lies in matching predictions with ground-truths, especially given the unordered nature of predicted results. We argue that the Hungarian algorithm [12], although widely adopted in the DETR families [5, 47, 28, 22, 13, 42], is not suitable for this task. Having a $O(n^3)$ time complexity and a $O(n^2)$ space complexity, the Hungarian algorithm is unable to handle a substantial number of voxels. In our experiments, associating two sets with 10K points each, the Hungarian algorithm consumes approximately 24 seconds and 2,304Mb of GPU memory on a 80G A100 GPU. In reality, the voxel number can go up to ∼70K in the Occ3D-nuScenes [38] dataset. Thus directly applying the Hungrian algorithm for set-to-set matching is infeasible in the occupancy prediction context.

But is accurate one-to-one association truly necessary for occupancy prediction? We recognize that the goal of one-to-one correspondence between prediction results and ground-truth annotation is to obtain supervision signals, essentially complete, precise *point locations*, and accurate *point classes*. The heavylifting of one-to-one association can be entirely avoided if we can obtain such supervision signals elsewhere. Therefore, we propose to decouple the occupancy prediction task into two parallel subtasks, as illustrated in Fig. 1. The first task obtains supervision on point locations by aligning predicted point distributions with ground-truths, a task achievable through the Chamfer distance loss, a well-established technique for point clouds [6, 30]. The second task obtains supervision on point classes by assigning semantic labels to predicted points. This is accomplished by assigning each point the class of its nearest neighbor in the ground-truths. It's noteworthy that all operations involved can be executed in parallel and are highly efficient on GPU devices. As a result, a single matching in Occ3D-nuScenes can be processed within milliseconds, with negligible memory consumption. With a time complexity of $O(n^2)$ and space complexity of $O(n)$, our formulation breaks the ground for large-scale training for the occupancy prediction models.

In addition, we propose several strategies to further boost the performance of occupancy prediction in our end-to-end sparse formulation, including coarse-to-fine learning, consistent point sampling, and adaptive loss re-weighting. On Occ3D-nuScenes, all our model variants easily surpass all prior work, verifying the efficacy and effectiveness of the proposed method. Especially, our most lightweight model achieves a 3.3 absolute RayIoU improvement compared with SparseOcc [21] while operating more than 2× faster. The heaviest configuration ultimately achieves a RayIoU of 41.2, establishing a new upper bound with a 14% advantage. Our contributions are summarized as follows:

- To the best of our knowledge, for the first time, we view the occupancy prediction as a direct set prediction problem, facilitating end-to-end training of the sparse framework.

2

- Several non-trivial strategies, including coarse-to-fine learning, consistent point sampling, and adaptive re-weighting, are further introduced for boosting the performance of OPUS.

- Extensive experiments on Occ3D-nuScenes reveal that OPUS can outperform state-of-the-art methods in terms of RayIoU results, while maintain a real-time inference speed.

## 2 Related work

### 2.1 Occupancy prediction

Occupancy prediction entails determining the occupancy status and class of each voxel within a 3D space. This task has recently become a foundational perception task in autonomous driving and raises great interests from both academic and industrial communities. Conventional methods [3, 46, 9, 27, 16, 41, 38, 4] typically employ the continuous and dense feature representation, which, however, suffer from computational redundancy due to the inherent sparsity of occupancy data. In addressing this issue, Tang et al. [37] compresses the dense feature using the Tri-Perspective View representation for model efficiency. Recently, several transformer-based approaches [21, 10, 14] with sparse queries have emerged. For example, OccupancyDETR [10] conducts object detection followed by assigning each object with one query for occupancy completion. VoxFormer [14] generates 3D voxels from a set of sparse queries, corresponding to occupied locations identified through a pre-task of depth estimation. Meanwhile, SparseOcc [21] employs a series of sparse voxel decoders to filter out empty grids and predict occupied statuses of retained voxels in each stage. While these approaches have succeeded in reducing computational costs, they often necessitate multi-stage processes and intricate space modeling. In contrast, our method directly applies sparse queries to regress the occupied locations without pre-defined locations, facilitating an elegant and end-to-end occupancy prediction.

### 2.2 Set prediction with transformers

The concept of directly predicting sets with Transformers was initially introduced by DETR [5], where a set of sparse queries generates unordered detection results with feature and object interactions. By viewing the object detection as a direct set prediction problem, DETR eliminates the need for complex post-processing, enabling end-to-end performance. Following DETR, numerous variants [47, 28, 22, 13, 42, 45, 36] have been proposed for performance improvements and efficient training. The effectiveness of the sparse-query-based paradigm has also been validated in 3D object detection [43, 23, 18–20, 40], where 3D information is encoded into the queries. For example, DETR3D [43] employs a sparse set of 3D object queries to index 2D features, linking 3D positions to multi-view images using camera transformation matrices. PETR [23] generates 3D position-aware features by encoding 3D position embedding into 2D image features, enabling queries to directly aggregate features without the 3D-to-2D projection. Sparse4D [18] further advances sparse 3D object detection by refining detection results with spatial-temporal feature fusion. Despite of the great success, set prediction with Transformers remains restricted primarily to object detection, where the query number are typically small due to the limited object number in a scene. Extending this approach to occupancy prediction poses a big challenge due to the substantially larger number of queries required.

## 3 Methodology

In this part, we first recap current query-based sparsification approaches for occupancy prediction in Sec. 3.1. Then, Sec. 3.2 describes our formulation that views the task as a direct set prediction problem. Finally, we detail the proposed OPUS framework in Sec. 3.3.

### 3.1 Revisiting query-based occupancy sparsification

Transformers with sparse queries offer a promising avenue for tackling the inherent sparsity in occupancy representation. A notable approach to reduce the number of queries is allocating each query to a patch of voxels rather than a single voxel, as presented in PETRv2 [24]. However, this method still generates a dense prediction of the 3D space, thus failing to efficiently address the redundancy issue. Alternatively, VoxFormer [14] and SparseOcc [21] allocate sparse queries exclusively to occupied voxels. VoxFormer employs a depth estimation module to identify potentially occupied voxels, while SparseOcc utilizes multiple stages to progressively filter out empty regions. Nonetheless, their sparsification processes rely on accurately recognizing the occupancy status of voxels and therefore suffer from the cumulative errors. Moreover, their pipelines necessitate intricate intermediate descriptions of the 3D space, hindering seamless end-to-end operation.

The dilemmas of current approaches significantly stem from treating the task as a classification problem, where each query is confined to a specific physical region for classifying the semantic labels. This constraint severely limits query flexibility, preventing adaptive focus on suitable areas. To address this, we propose to remove this restriction by allowing each query to autonomously determine its relevant area. In the end, we view occupancy prediction as a direct set prediction problem, where each query predicts point positions and semantic classes, simultaneously.

## 3.2 A set prediction problem

At the core of our work lies the conceptualization of occupancy prediction as a set prediction task. We denote the $V_g$ occupied voxels in the ground-truth as $\{\mathbb{P}_g, \mathbb{C}_g\}$, where $|\mathbb{P}_g| = |\mathbb{C}_g| = V_g$. For each entry in $\{\mathbf{p}_g, \mathbf{c}_g\} \in \{\mathbb{P}_g, \mathbb{C}_g\}$, $\mathbf{p}_g$ represents the 3D coordinates of a voxel center, while $\mathbf{c}_g$ stores the semantic class of the corresponding voxel. Given the predictions $\{\mathbb{P}, \mathbb{C}\}$ of $V$ points, our primary challenge is to devise an effective strategy for set-to-set matching. In other words, we must determine how to supervise the training of unordered predictions with the ground-truth data. One alternative is to adopt the Hungarian algorithm. However, our previous discussions and experiments in the appendix reveal its scalability limitations. Rather than pursuing one-to-one associations between the predictions and ground-truths, we recognize the matching essentially aims at accurate locations and classes in predictions. This motivates us to decouple the task into two parallel objectives: (1) Encouraging the predicted locations to be precise and comprehensive. (2) Ensuring the predicted points are assigned with proper semantic classes from the ground-truth labels.

The first objective focuses on aligning distributions between predicted and ground-truth points, a task achievable through the Chamfer distance loss which is well-proved in the field of point clouds [6, 30, 11, 44]:

$$\text{CD}(\mathbb{P}, \mathbb{P}_g) = \frac{1}{|\mathbb{P}|} \sum_{\mathbf{p} \in \mathbb{P}} D(\mathbf{p}, \mathbb{P}_g) + \frac{1}{|\mathbb{P}_g|} \sum_{\mathbf{p}_g \in \mathbb{P}_g} D(\mathbf{p}_g, \mathbb{P}), \text{ where } D(\mathbf{x}, \mathbb{Y}) = \min_{\mathbf{y} \in \mathbb{Y}} ||\mathbf{x} - \mathbf{y}||_1. \quad (1)$$

Minimizing Chamfer distance leads to similar distributions of predictions and ground-truths, enabling direct learning of occupied voxels without necessitating knowledge of their orders.

Concerning the second objective, although direct comparison between $\mathbb{C}$ and $\mathbb{C}_g$ is invalid due to their correspondence to different locations, we can leverage the spatial locality properties of voxels to find a proxy. Nearby points belonging to the same object usually carry the same semantic labels, thus we propose assigning each predicted point the class of its nearest neighbor voxel in the ground-truth:

$$\{\hat{\mathbb{C}}, \hat{\mathbb{P}}\} = \left\{ \arg\min_{\{\mathbf{c}_g, \mathbf{p}_g\} \in \{\mathbb{C}_g, \mathbb{P}_g\}} ||\mathbf{p}_g - \mathbf{p}||_2, \quad \mathbf{p} \in \mathbb{P} \right\}. \quad (2)$$

Here, $\hat{\mathbb{C}}$ is the updated classes that are prepared to supervise the learning of the predicted $\mathbb{C}$.

It's noteworthy that computations of both Eq. (1) and Eq. (2) can be executed efficiently and in parallel on GPU devices. As a result, a single matching can be swiftly processed within milliseconds, enabling feasibility of the large-scale training by viewing the occupancy prediction task as a direct set prediction problem. Next, we delve into the specifics of the proposed OPUS framework.

## 3.3 Details of OPUS

This part describes OPUS framework, as illustrated in Fig. 2. Initially, image features are extracted from multi-view images. And a set of learnable queries $\mathbb{Q}$, point positions $\mathbb{P}$, and scores $\mathbb{C}$ are initialized. Subsequently, these query features and prediction outcomes are fed into a sequence of decoders, undergoing iterative refinement through correlation with image features. At each stage, predicted positions and scores are supervised by the ground-truths, facilitating end-to-end training for the entire framework. It can be observed that our most important structure is the sequence of multiple decoders. Therefore, we next provide a detailed description to the inputs/outputs of the decoders and how features are aggregated and updated within the decoders.

**Notations.** Denote the set of learnable queries, point positions, and point scores as $\{\mathbb{Q}_0, \mathbb{P}_0, \mathbb{C}_0\}$ before feeding into decoders, and as $\{\mathbb{Q}_i, \mathbb{P}_i, \mathbb{C}_i\}$ for the outcomes of the $i$-th decoder. The length of these sets is all $Q$, which corresponds to the number of queries. Each query feature $\mathbf{q}_i \in \mathbb{Q}_i, i \in \{0, 1, \cdots, 6\}$ has a channel size $C$, set to 256 in our implementation. To reduce the number of queries, which is a bottleneck for model efficiency, each query $q_i$ predicts $R_i$ points rather than a single one. Consequently, $\mathbf{p}_i \in \mathbb{P}_i$ and $\mathbf{c}_i \in \mathbb{C}_i$ have shapes of $Q \times R_i \times 3$ and $Q \times R_i \times N$, respectively. Here, $N$ represents the number of semantic classes.
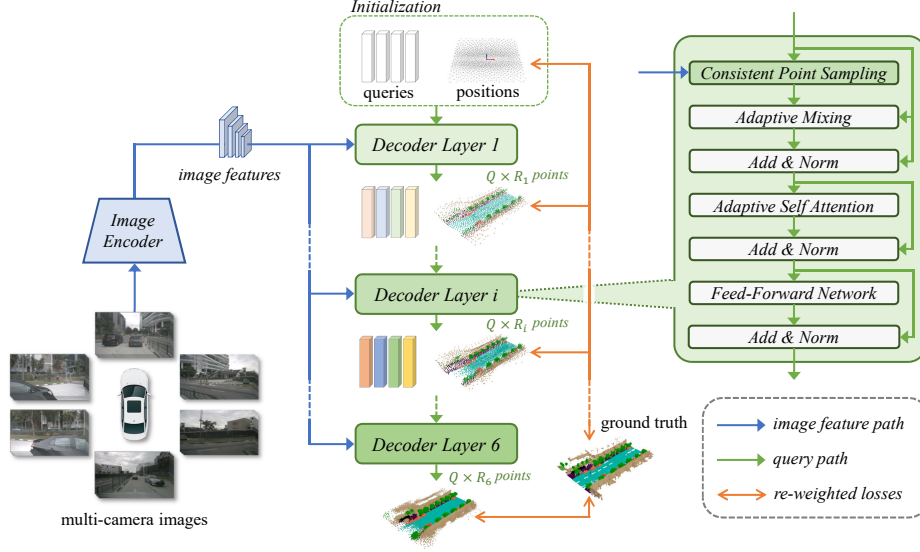
4

Figure 2: OPUS leverages a transformer encoder-decoder architecture comprising: (1) An image encoder to extract 2D features from multi-view images. (2) A series of decoders to refine the queries with image features, which are correlated via the *consistent point sampling* module. (3) A set of learnable queries to predict locations and classes of occupancy points. Each query obeys a *coarse-to-fine* rule, progressively increasing the number of predicted points. In the end, the entire model is trained end-to-end using our *adaptively re-weighted* set-to-set losses.

**Coarse-to-fine prediction.** High-level semantic information can be difficult to predict accurately from just low-level features. Therefore, instead of attempting to predict occupancy for the entire 3D environment, we allow the model to predict "sparse" occupancy results in early stages, as shown in Fig. 2. To achieve this, we follow a coarse-to-fine strategy, gradually increasing the number of points generated from one query. In other words, we always have $R_{i-1} \leq R_i$ for $i \in \{1, 2, \cdots, 6\}$.

It's noteworthy that the Chamfer distance has another advantage over the Hungarian algorithm here: even when the number of predictions is smaller than that of the ground-truths, the assignment won't collapse into a local shape of the ground-truths. This is because the Hungarian algorithm could assign the predictions to any subset of the ground-truths due to its lack of distribution constraints. In contrast, the Chamfer distance maintains a global perspective, considering the overall distribution of points rather than enforcing a strict one-to-one correspondence. This ensures that the predicted points are more evenly distributed and representative of the actual 3D environment, even when fewer in number.

**Details of the decoder.** Our decoder is analogous to that in SparseBEV [20], a performant and sparse object detector. For a given query $\mathbf{q}_{i-1} \in \mathbb{Q}_{i-1}$ and its corresponding point locations $\mathbf{p}_{i-1} \in \mathbb{P}_{i-1}$, the $i$-th decoder first aggregates image features through a consistent point sampling, a new scheme elaborated in our subsequent discussion. Subsequently, the query feature is updated into $\mathbf{q}_i$ with the adaptive mixing of image and query features, along with the self-attention among all queries, mirroring operations in SparseBEV. In the end, a prediction module, comprising only Linear, LayerNorm, and ReLU layers, generates the semantic classes $\mathbf{c}_i$ (size $R_i \times N$) and the position offsets $\Delta\mathbf{p}_i$ (size $R_i \times 3$). As the $\Delta\mathbf{p}_i$ cannot be directly added to $p_{i-1}$ due to dimension misalignment, we first compute the mean of $\mathbf{p}_{i-1}$ along the first dimension and then duplicate the results by $R_i$ times into $\bar{\mathbf{p}}_{i-1}$. The final position $\mathbf{p}_i$ is computed as $\mathbf{p}_i = \bar{\mathbf{p}}_{i-1} + \Delta\mathbf{p}_i$.

**Consistent point sampling.** The feature sampling method utilized in SparseBEV is not applicable for our method as it is specifically designed for detection inputs. Therefore, we propose a novel process of Consistent Point Sampling (CPS), aiming at sampling 3D points and aggregating features from $M$ image features. Given input $\{\mathbf{q}, \mathbf{p}\} \in \{\mathbb{Q}, \mathbb{P}\}$, we sample $S$ points and find their respective coordinates in the $m$-th image feature by the following equation:

$$\mathbf{c}_m = \mathbf{T_m r}, \text{ where } \mathbf{r} = \mathbf{m}_p + \phi(\mathbf{q}) \cdot \sigma_p, \tag{3}$$

where $\mathbf{T_m}$ represents the projection matrix from current 3D space into the $m$-th image's coordinates. $\phi(\mathbf{q})$ generates $S$ 3D points from the query feature $\mathbf{q}$ using a linear layer. $\mathbf{m}_p$ and $\sigma_p$ denote the mean

5

and standard deviation, respectively, of the $R$ points in $\mathbf{p}$. It's worthy to note that we re-weight the predicted offsets $\phi(\mathbf{q})$ with the standard deviation $\sigma_p$ to inherent the dispersion degree from previous predictions. In essence, we tend to sample more aggressively if the input $\mathbf{p}$ contains diverse points, and sample points in a narrower range otherwise. This operation can evidently enhance the prediction performance, as demonstrated in our experiments.

Not all coordinates in $\mathbf{c}_m$ are feasible since the sampled points might not be visible within the corresponding camera. Therefore, we generate a mask set $\mathbb{V}_m$ where the $s$-th value is 1 if $\mathbf{c}_{s,m}$ is valid and 0 otherwise, for $s \in \{1, 2, \cdots, S\}$ and $m \in \{1, 2, \cdots, M\}$. Next, we aggregate information from image features $\{F_m\}_1^M$ for the later adaptively mixing stage. Specifically, we have

$$f_s = \frac{1}{\sum_{m=1}^{M} |\mathbb{V}_m|} \sum_{s=1}^{S} \sum_{m=1}^{M} w_{s,m} \cdot v_{s,m} \cdot \mathcal{B}(F_m, \mathbf{c}_{s,m}), \tag{4}$$

where $v_{s,m}$ denotes the $s$-th element in $\mathbb{V}_m$ and $\mathbf{c}_{s,m}$ is the coordinates of the $s$-th point $\mathbf{r}_s$ mapped into the $m$-th image feature. The operation $\mathcal{B}$ refers to the bilinear interpolation. $w_{s,m}$ is the weight for the $\mathbf{r}_s$ on the $m$-th image feature, generated from the query feature $\mathbf{q}$ by linear transformation.

**The training loss with adaptively re-weighting.** The training object of our framework is to supervise the learning of $\{\mathbb{P}_i, \mathbb{C}_i\}_{i=1}^{6}$ with the ground-truth $\{\mathbb{P}_g, \mathbb{C}_g\}$. Point positions can be trained with Eq. (1). However, the original Chamfer distance loss focuses on the overall similarity of point distributions, neglecting whether each individual is good enough. This leads to unsatisfactory performance, as observed in our experiments. To cope with this issue, we employ a simple but effective re-weighting strategy to emphasize erroneous points, and modify the Chamfer distance loss as follows:

$$\mathrm{CD}_R(\mathbb{P}, \mathbb{P}_g) = \frac{1}{|\mathbb{P}|} \sum_{\mathbf{p} \in \mathbb{P}} D_R(\mathbf{p}, \mathbb{P}_g) + \frac{1}{|\mathbb{P}_g|} \sum_{\mathbf{p}_g \in \mathbb{P}_g} D_R(\mathbf{p}_g, \mathbb{P}),$$
$$\text{where } D_R(\mathbf{x}, \mathbb{Y}) = W(d) \cdot d \text{ with } d = \min_{\mathbf{y} \in \mathbb{Y}} ||\mathbf{x} - \mathbf{y}||_1. \tag{5}$$

Here, $W(d)$ is the re-weighting function penalizing points with large distance to the closest ground-truths. In our implementation, we use a step function of $W(d)$ being 5 if $d \geq 0.2$ and 1 otherwise.

For the classification, we first generate the target classes $\hat{\mathbb{C}}_i$ for $\mathbb{C}_i$ using Eq. (2). Subsequently, the semantic classes can be trained with the conventional classification losses. In our implementation, we adopt the focal loss [17] with mannually searched weights on different categories and denote the modified loss as $\mathrm{FocalLoss}_R$. In the end, the training objective of the proposed OPUS becomes

$$L_{\mathrm{OPUS}} = \mathrm{CD}_R(\mathbb{P}_0, \mathbb{P}_g) + \sum_{i=1}^{6} (\mathrm{CD}_R(\mathbb{P}_i, \mathbb{P}_g) + \mathrm{FocalLoss}_R(\mathbb{C}_i, \hat{\mathbb{C}}_i)), \tag{6}$$

where $\mathrm{CD}_R(\mathbb{P}_0, \mathbb{P}_g)$ explicitly encourages initial points $\mathbb{P}_0$ to capture a general pattern of the dataset.

## 4 Experiments

### 4.1 Experimental setup

**Dataset and metrics.** All models are evaluated on the Occ3D-nuScenes [38] dataset, which provides occupancy labels for 18 classes (1 free class and 17 semantic classes) on the large-scale nuScenes [2] benchmark. Out of the 1,000 labeled driving scenes, 750/150/150 are used for training/validation/testing, respectively. The commonly used mIoU metric is utilized for evaluation. Recently, SparseOcc [21] points that that overestimation can easily hack the mIoU metric and proposes RayIoU as a remedy. Therefore, following their work, we also report the RayIoU results under different distance thresholds at 1, 2, and 4 meters, denoted as $\mathrm{RayIoU}_{1m}$, $\mathrm{RayIoU}_{2m}$, and $\mathrm{RayIoU}_{4m}$, respectively. The final RayIoU score is the average of these three values.

**Implementation details.** Following previous works [21, 16, 8], we resize images to $704 \times 256$ and extract features using a ResNet50 [7] backbone. We denote a series of models as OPUS-T, OPUS-S, OPUS-M and OPUS-L, with 0.6K, 1.2K, 2.4K and 4.8K queries, respectively. In each model, all queries predict an equal number of points, totalling 76.8K points in the final stage. The sampling number in our CPS is 4 for OPUS-T and 2 for other models. Please refer to Appendix D.2 for more details of different models. All models are trained on 8 nvidia 4090 GPUs with a batch size of 8 using the AdamW [26] optimizer. The learning rate warms up to $2e^{-4}$ in the first 500 iterations and then decays with a Cosine Annealing [25] scheme. Unless otherwise stated, models in main results are trained for 100 epochs and those in the ablation study are trained for 12 epochs.

Table 1: Occupancy prediction performance on Occ3D-nuScenes [38]. "8f" and "16f" denote models fusing temporal information from 8 or 16 frames, respectively. Baseline results are directly copied from their corresponding papers or the SparseOcc [21]. FPS results are measured on an A100 GPU.

| Methods | Backbone | Image Size | mIoU | RayIoU$_{1m}$ | RayIoU$_{2m}$ | RayIoU$_{4m}$ | RayIoU | FPS |
|---------|----------|-----------|------|---------|---------|---------|--------|-----|
| RenderOcc [31] | Swin-B | $1408 \times 512$ | 24.5 | 13.4 | 19.6 | 25.5 | 19.5 | - |
| BEVFormer [15] | R101 | $1600 \times 900$ | 39.3 | 26.1 | 32.9 | 38.0 | 32.4 | 3.0 |
| BEVDet-Occ [8] | R50 | $704 \times 256$ | 36.1 | 23.6 | 30.0 | 35.1 | 29.6 | 2.6 |
| BEVDet-Occ (8f) [8] | R50 | $704 \times 384$ | 39.3 | 26.6 | 33.1 | 38.2 | 32.6 | 0.8 |
| FB-Occ (16f) [16] | R50 | $704 \times 256$ | 39.1 | 26.7 | 34.1 | 39.7 | 33.5 | 10.3 |
| SparseOcc (8f) [21] | R50 | $704 \times 256$ | - | 28.0 | 34.7 | 39.4 | 34.0 | 17.3 |
| SparseOcc (16f) [21] | R50 | $704 \times 256$ | 30.6 | 29.1 | 35.8 | 40.3 | 35.1 | 12.5 |
| OPUS-T (8f) | R50 | $704 \times 256$ | 33.2 | 31.7 | 39.2 | 44.3 | 38.4 | 22.4 |
| OPUS-S (8f) | R50 | $704 \times 256$ | 34.2 | 32.6 | 39.9 | 44.7 | 39.1 | 20.7 |
| OPUS-M (8f) | R50 | $704 \times 256$ | 35.6 | 33.7 | 41.1 | 46.0 | 40.3 | 13.4 |
| OPUS-L (8f) | R50 | $704 \times 256$ | 36.2 | 34.7 | 42.1 | 46.7 | 41.2 | 7.2 |



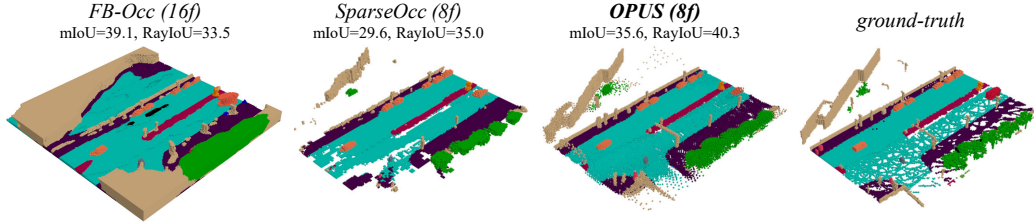| *FB-Occ (16f)* | *SparseOcc (8f)* | **OPUS (8f)** | *ground-truth* |
|---|---|---|---|
| mIoU=39.1, RayIoU=33.5 | mIoU=29.6, RayIoU=35.0 | mIoU=35.6, RayIoU=40.3 | |

Figure 3: Visualizations of occupancy predictions. **Best viewed in color**.

## 4.2 Main results

**Quantitative Performances.** In this part, we compare OPUS with previous state-of-the-art methods on the Occ3D-nuScenes dataset. Our methods not only achieves the superior performances in terms of RayIoU and competitive results in mIoU, but also demonstrates commendable real-time performance. As depicted in Tab. 1, OPUS-T (8f) reaches 22.4 FPS, significantly faster than dense counterparts and nearly 1.3 times the speed of sparse counterpart SparseOcc (8f). Despite using only 7 history frames, its 38.4 RayIoU result easily outperforms other models, including FB-Occ (16f) with RayIoU of $33.5(-4.9)$ and SparseOcc (16f) with RayIoU of $35.1(-3.3)$. Similarly, OPUS-S (8f) and OPUS-M (8f) achieve a good balance between performance and efficiency. The heaviest version of OPUS ultimately achieves an RayIoU of 41.2, surpassing the previous best result by a notable margin of 6.1.

With the same total number of points predicted, we vary the query number and correspondingly change the number of points from each query, leading to different versions of OPUS. It can be observed that increasing the query number decreases the FPS values from 22.4 to 7.2, while simultaneously boosts model performance in terms of mIoU and RayIoU. The OPUS-M (8f), with 2.4K queries, strikes a balance by achieving a comparable RayIoU while maintaining competitive FPS.

Despite the vulnerability of mIoU metric to overestimation manipulations [21], our OPUS attains a comparable mIoU of 36.2, significantly bridging the gap between dense and sparse models in this metric. These results under different metrics collectively demonstrate the superiority of our OPUS.

**Visualization.** We visualize the predicted occupancy in Fig. 3. It can be observed that FB-Occ tends to produce denser results compared to sparse methods. Though seems complete in the 3D environment, its predicted occupancy results are severely over-estimated, especially for the far areas. The overestimation may hack the mIoU metric [21], while heavily penalized by RayIoU that primarily considers the first occupied voxels along rays. Consequently, FB-Occ achieves the best mIoU of 39.1 but the worst RayIoU value. On the other hand, SparseOcc occasionally exhibits discontinuous predictions with false negatives, especially in long distances. This is attributed to SparseOcc's gradual removal of empty voxels, making erroneous filtering in early stages accumulates and contributes to the final false predictions. In contrast, our OPUS maintains a more continuous prediction thanks to its end-to-end approach, resulting in a more reasonable visualization.

## 4.3 Ablation study and visualizations

This part details our ablation study and visualizations using the OPUS-M (8f) model.

Table 2: Model performances with different combinations of proposed strategies.

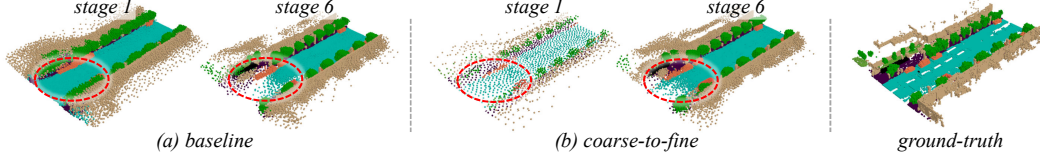| $CD_R$ | $FocalLoss_R$ | CPS | Coarse-to-fine | mIoU | RayIoU$_{1m}$ | RayIoU$_{2m}$ | RayIoU$_{4m}$ | RayIoU |
|---|---|---|---|---|---|---|---|---|
| | | | | 17.4 | 23.6 | 29.7 | 34.3 | 29.2 |
| ✓ | | | | 23.7 (6.3↑) | 23.9 | 30.7 | 35.6 | 30.1 (0.9↑) |
| ✓ | ✓ | | | 25.1 (1.4↑) | 25.2 | 32.3 | 37.0 | 31.5 (1.4↑) |
| ✓ | ✓ | ✓ | | 25.5 (0.4↑) | 26.0 | 33.1 | 37.9 | 32.3 (0.8↑) |
| ✓ | ✓ | ✓ | ✓ | 27.2 (1.7↑) | 26.1 | 33.3 | 38.4 | 32.6 (0.3↑) |



Figure 4: Visualizations of the coarse-to-fine predictions.

**Effects of the proposed strategies in OPUS.** In our work, we introduce adaptive re-weighting for the Chamfer distance loss and focal loss, along with consistent point sampling, and coarse-to-fine prediction strategies. We examine the impacts of these strategies as shown in Tab. 2. Without bells and whistles, OPUS achieves a baseline 17.4 mIoU and a 29.2 RayIoU. Replacing the original CD loss into our revision $CD_R$ significantly boosts the mIoU and RayIoU by 6.4 and 0.9, respectively, demonstrating the importance of focusing on erroneous predicted locations in this task. The $FocalLoss_R$ further improves both metrics by 1.4. Incorporating the term $\sigma_p$ in Eq. (3) further enhances mIoU and RayIoU by 0.4 and 0.8, demonstrating the efficacy of considering previous point distribution in the current sampling process. The proposed coarse-to-fine query prediction gradually increases the number of points across the stages. This scheme not only reduces computations in early stages but also notably benefits model performance, particularly in mIoU, which is increased by 1.7. These results highlight the cumulative benefits of each component, showcasing how their integration leads to substantial performance gains.

**Visualization on the coarse-to-fine prediction.** We visualize the prediction results at different stages in Fig. 4. In the baseline scenario depicted in Fig. 4(a), where all decoders regress the same number of points, we observe inconsistent point distributions across stages and numerous false negative predictions in long distances, as highlighted by circles. This may be attributed to the difficulty of learning the fine-grained occupancy representations in the early stages, impeding the efficient training of the entire framework. In contrast, our coarse-to-fine strategy significantly alleviates the learning difficulty in early stages, thereby leading to improved model performances. As a result, the point distributions are more consistent among different stages, and the final predictions exhibit much fewer false negatives, as illustrated in Fig. 4(b).

**Visualizations of predicted points.** In Fig. 6, we select a few queries and visualize their predicted points. Notably, most queries exhibit a tendency to predict points with consistent classes, or even from the same instance, as depicted in Fig. 6(a)-(g). An interesting observation is that the predicted points tend to exhibit diverse distributions in classes with large volumes, such as drivable surfaces and sidewalks. Conversely, for objects with limited sizes, such as traffic cones, motorcycles, and cars, the points are distributed more closely with respect to the instance size. The patterns can be further verified by Fig. 5, where we present the standard deviations of points from queries with three chosen classes. These results highlight the efficacy of our model in adapting its predictions to the distinct spatial characteristics of various object classes.
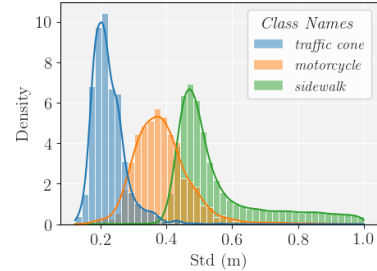


Figure 5: Distributions of standard deviations of points from one query.

As we do not explicitly constrain points from one query to have the same class, it's conceivable that one query could yield points of different classes. We found this phenomenon commonly occurs at the boundaries between objects. However, even when classes vary, these points are still closely distributed, as depicted in Fig. 6(h)-(j).

Table 3: Comparison of various treatments on initial locations $\mathbb{P}^0$. "Grid" and "Random" indicate that points are sampled uniformly in BEV space and randomly in the 3D space, respectively "Optimized" means that points are randomly initialized but supervised with ground-truths via the $\text{CD}_R$ loss.

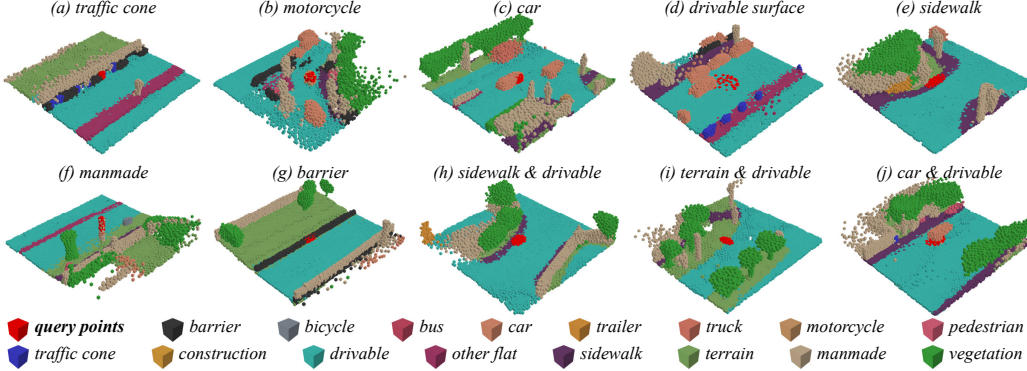| Type | mIoU | RayIoU$_{1m}$ | RayIoU$_{2m}$ | RayIoU$_{4m}$ | RayIoU |
|------|------|------|------|------|------|
| Grid | 22.8 | 22.2 | 28.9 | 33.9 | 28.3 |
| Random | 23.1 | 23.6 | 30.5 | 35.6 | 29.9 |
| Optimized | 23.7 | 23.9 | 30.7 | 35.6 | 30.1 |



Figure 6: Visualizations of points generated from different queries. **Best viewed in color**.

**Influence of treatments on the initial points.** Tab. 3 compares three different treatments on the initial points $\mathbb{P}^0$. Grid initialization divides the BEV space into evenly-distributed pillars and orderly assigns pillar centers as the initial locations, a method utilized in BEVFormer [15]. Random initialization assigns each location with a uniform distribution in the 3D space. After initialization, $\mathbb{P}^0$ remains learnable during training. On top of the random initialization, our OPUS further add supervisions of the ground-truth distributions to $\mathbb{P}^0$ (*i.e.*, $\text{CD}_R(\mathbb{P}^0, \mathbb{P}_g)$ in Eq. (6)). The results in Tab. 4 show that random initialization outperforms grid initialization, achieving an mIoU of 23.1 compared to 22.8, and a RayIoU of 29.9 compared to 28.3. This improvement is likely due to the fact the random initialization provides a more diverse 3D distribution. Furthermore, the introduced supervision results in additional improvements of 0.6 on mIoU and 0.2 on RayIoU. These results reveal the efficiency of the random initialization and the additional supervision on the initial locations.

**Visualization of the self-attention.** For better understanding what query points attend to, we visualize top 10 query points with highest self attention weights for each query. We project them on 2D image for better visualization. Here are some interesting finding from Fig. 7. Generally speaking, the query tends to attend to neighbouring query points. For example, the sidewalk query in the first image and the car query in the second image, allow local information from neighbouring query points flow into, enabling the query to capture detailed local information. Additionally, the query maintains the ability to attend to semantic related locations even if they are not very close to the query point. For instance, the vegetation query in the first image of Fig. 7 not only attends to the stem of the tree, but also the grass, indicating that the query can capture semantically related information for more accurate predictions. Another notable observation is that query can sense terrain-related information. For instance, in the third image, the sidewalk query points attend to other points along a straight line following the road's edge, highlighting the model's ability to understand scene related structure of the environment.



Figure 7: the self-attention in decoders. for each pivot (marked as $\times$), query points with top 10 attention weights are shown by circles, with sizes proportional to weights. **best viewed in color**.

Table 4: Comparisons between different sparsification strategies.

| Model | $Q$ | $R$ | RayIoU$_{1m}$ | RayIoU$_{2m}$ | RayIoU$_{4m}$ | RayIoU | FPS |
|-------|-----|-----|---------------|---------------|---------------|--------|-----|
| SparseOcc | (4000/16000/64000) | | 28.4 | 34.9 | 39.6 | 34.3 | 17.3 |
| PETR v2 | 2500 | 256 | 24.4 | 31.0 | 36.3 | 30.6 | 13.8 |
| OPUS | 2400 | 32 | 31.7 | 38.8 | 43.4 | 38.0 | 13.4 |

Table 5: Performance on the Waymo-Occ3D dataset.

| Model | General | Vehicle | Bicyclist | Ped. | Sign | Tfc. light | Pole | Cons. cone | Bicycle | Motorcycle | Building | Vegetation | Treetrunk | Road | Sidewalk | mIoU | RayIoU | FPS |
|-------|---------|---------|-----------|------|------|------------|------|------------|---------|------------|----------|------------|-----------|------|----------|------|--------|-----|
| BEVDet | 0.13 | 13.06 | 2.17 | 10.15 | 7.80 | 5.85 | 4.62 | 0.94 | 1.49 | 0.0 | 7.27 | 10.06 | 2.35 | 48.15 | 34.12 | 9.88 | - | - |
| TPVFormer | 3.89 | 17.86 | 12.03 | 5.67 | 13.64 | 8.49 | 8.90 | 9.95 | 14.79 | 0.32 | 13.82 | 11.44 | 5.8 | 73.3 | 51.49 | 16.76 | - | - |
| BEVFormer | 3.48 | 17.18 | 13.87 | 5.9 | 13.84 | 2.7 | 9.82 | 12.2 | 13.99 | 0.0 | 13.38 | 11.66 | 6.73 | 74.97 | 51.61 | 16.76 | - | 4.6 |
| CTF-Occ | 6.26 | 28.09 | 14.66 | 8.22 | 15.44 | 10.53 | 11.78 | 13.62 | 16.45 | 0.65 | 18.63 | 17.3 | 8.29 | 67.99 | 42.98 | 18.73 | - | 2.6 |
| OPUS-L | 4.66 | 27.07 | 19.39 | 6.53 | 18.66 | 6.41 | 11.44 | 10.40 | 12.90 | 0.0 | 18.73 | 18.11 | 7.46 | 72.86 | 50.31 | 19.00 | 24.7 | 8.5 |

**Comparisons between different sparsification strategies** In Tab. 4, we compare OPUS to two other models with different sparsification strategies. The first baseline is SparseOcc, which achieves sparsification by filtering out empty voxels at various cascade stages. Following PETRv2 [24], the second baseline is a pillar-patch based method that partitions the 3D space into a small number of pillar-patches. We use $50 \times 50$ queries with each corresponding to the classification of neighbouring $4 \times 4 \times 16$ voxels. For a fair comparison, all these models are trained for 100 epochs. In contrast, our model achieve best results after sufficient training with a RayIoU score of 38.0,far outperforming SparseOcc with a RayIoU score of 34.3. On the other hand, our model can also runs in a real-time speed. These results demonstrates the superiority of our sparification procedure.

**Comparisons on the Waymo-Occ3D dataset.** We further simply implement OPUS on the Waymo-Occ3D [35] dataset to explore the generalization and robustness of OPUS. As Waymo-Occ3D is not commonly used as a standard benchmark for vision-centric approaches, the only vision-based method we found with reported results on this dataset is the Occ3D paper, which evaluates BEVDet, TPVFormer, BEVFormer, and the newly proposed CTF-Occ [38]. We trained the OPUS-L (1f) on 20% of the dataset for a fair comparison with these baselines. As reported in Tab. 5, despite not fine-tuning the training configurations, OPUS-L already achieves 19.0 mIoU, outperforming all previous methods. Moreover, OPUS-L also reaches 8.5 FPS on the Waymo-Occ3D dataset, which is around 3 times the speed of CTF-Occ and 2 times the speed of BEVFormer.

## 5 Conclusions and limitations

This paper introduces a novel perspective on occupancy prediction by framing it as a direct set prediction problem. Using a transformer encoder-decoder architecture, the proposed OPUS directly predicts occupied locations and classes in parallel from a set of learnable queries. The matching between predictions and ground truths is accomplished through two efficient tasks in parallel, facilitating end-to-end training with a large number of points in this application. In addition, the query features are enhanced via a list of non-trivial designs (*i.e.*, coarse-to-fine learning, consistent point sampling, and loss re-weighting), and therefore leads to boosted prediction performances. Our experiments on the Occ3D-nuScenes benchmark demonstrate that OPUS surpasses all prior arts in terms of both accuracy and efficiency, thanks to the sparse designs in our framework.

However, the proposed OPUS also comes with new challenges, particularly regarding the convergence speed. The slow convergence may potentially be alleviated by drawing lessons from follow-up works of DETR, which have largely addressed the convergence issue of the original DETR. Another challenge is that while sparse approaches typically achieve higher RayIoU compared to dense counterparts, they often struggle with the mIoU metric. Improving the mIoU performance while maintaining superior RayIoU results is a promising direction for future works. Moreover, despite conducting experiments on vision-only datasets, our core formulation is directly applicable to multi-modal tasks as well. We leave the multi-modal occupancy prediction as future work.

## Acknowledgments

## 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.

[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.

[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*, pages 3991–4001, 2022.

[4] Qinglong Cao, Yuntian Chen, Ding Wang, Zhengqin Xu, Chao Ma, Xiaokang Yang, and Shiyi Chen. Vision-informed flow field super-resolution with quaternion spatial modeling and dynamic fluid convolution. *Physics of Fluids*, 36(9), 2024.

[5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020.

[6] Haoqiang Fan, Hao Su, and Leonidas J Guibas. A point set generation network for 3d object reconstruction from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 605–613, 2017.

[7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

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

[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.

[10] Yupeng Jia, Jie He, Runze Chen, Fang Zhao, and Haiyong Luo. Occupancydetr: Making semantic scene completion as straightforward as object detection. *arXiv preprint arXiv:2309.08504*, 2023.

[11] Tarasha Khurana, Peiyun Hu, David Held, and Deva Ramanan. Point cloud forecasting as a proxy for 4d occupancy forecasting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1116–1124, 2023.

[12] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2 (1-2):83–97, 1955.

[13] Feng Li, Hao Zhang, Shilong Liu, Jian Guo, Lionel M Ni, and Lei Zhang. Dn-detr: Accelerate detr training by introducing query denoising. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13619–13627, 2022.

[14] 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*, pages 9087–9098, 2023.

[15] 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.

[16] Zhiqi Li, Zhiding Yu, David Austin, Mingsheng Fang, Shiyi Lan, Jan Kautz, and Jose M Alvarez. Fb-occ: 3d occupancy prediction based on forward-backward view transformation. *arXiv preprint arXiv:2307.01492*, 2023.

[17] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.

[18] 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.

[19] Xuewu Lin, Tianwei Lin, Zixiang Pei, Lichao Huang, and Zhizhong Su. Sparse4d v2: Recurrent temporal fusion with sparse model. *arXiv preprint arXiv:2305.14018*, 2023.

[20] Haisong Liu, Yao Teng, Tao Lu, Haiguang Wang, and Limin Wang. Sparsebev: High-performance sparse 3d object detection from multi-camera videos. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 18580–18590, 2023.

[21] Haisong Liu, Haiguang Wang, Yang Chen, Zetong Yang, Jia Zeng, Li Chen, and Limin Wang. Fully sparse 3d panoptic occupancy prediction. *arXiv preprint arXiv:2312.17118*, 2023.

[22] Shilong Liu, Feng Li, Hao Zhang, Xiao Yang, Xianbiao Qi, Hang Su, Jun Zhu, and Lei Zhang. Dab-detr: Dynamic anchor boxes are better queries for detr. *arXiv preprint arXiv:2201.12329*, 2022.

[23] 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.

[24] 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.

[25] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.

[26] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018.

[27] Qihang Ma, Xin Tan, Yanyun Qu, Lizhuang Ma, Zhizhong Zhang, and Yuan Xie. Cotr: Compact occupancy transformer for vision-based 3d occupancy prediction. *arXiv preprint arXiv:2312.01919*, 2023.

[28] Depu Meng, Xiaokang Chen, Zejia Fan, Gang Zeng, Houqiang Li, Yuhui Yuan, Lei Sun, and Jingdong Wang. Conditional detr for fast training convergence. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3651–3660, 2021.

[29] Qiang Meng, Xiao Wang, JiaBao Wang, Liujiang Yan, and Ke Wang. Small, versatile and mighty: A range-view perception framework. *arXiv preprint arXiv:2403.00325*, 2024.

[30] Benedikt Mersch, Xieyuanli Chen, Jens Behley, and Cyrill Stachniss. Self-supervised point cloud prediction using 3d spatio-temporal convolutional networks. In *Conference on Robot Learning*, pages 1444–1454. PMLR, 2022.

[31] 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. *arXiv preprint arXiv:2309.09502*, 2023.

[32] Yining Shi, Jingyan Shen, Yifan Sun, Yunlong Wang, Jiaxin Li, Shiqi Sun, Kun Jiang, and Diange Yang. Srcn3d: Sparse r-cnn 3d for compact convolutional multi-view 3d object detection and tracking. *arXiv preprint arXiv:2206.14451*, 2022.

[33] Yining Shi, Kun Jiang, Ke Wang, Kangan Qian, Yunlong Wang, Jiusi Li, Tuopu Wen, Mengmeng Yang, Yiliang Xu, and Diange Yang. Effocc: A minimal baseline for efficient fusion-based 3d occupancy network. *arXiv preprint arXiv:2406.07042*, 2024.

[34] Yining Shi, Jiusi Li, Kun Jiang, Ke Wang, Yunlong Wang, Mengmeng Yang, and Diange Yang. Panossc: Exploring monocular panoptic 3d scene reconstruction for autonomous driving. In *2024 International Conference on 3D Vision (3DV)*, pages 1219–1228. IEEE, 2024.

[35] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2446–2454, 2020.

[36] Zhiqing Sun, Shengcao Cao, Yiming Yang, and Kris M Kitani. Rethinking transformer-based set prediction for object detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3611–3620, 2021.

[37] Pin Tang, Zhongdao Wang, Guoqing Wang, Jilai Zheng, Xiangxuan Ren, Bailan Feng, and Chao Ma. Sparseocc: Rethinking sparse latent representation for vision-based semantic occupancy prediction. *arXiv preprint arXiv:2404.09502*, 2024.

[38] Xiaoyu Tian, Tao Jiang, Longfei Yun, Yucheng Mao, Huitong Yang, Yue Wang, Yilun Wang, and Hang Zhao. Occ3d: A large-scale 3d occupancy prediction benchmark for autonomous driving. *Advances in Neural Information Processing Systems*, 36, 2024.

[39] Jiabao Wang, Qiang Meng, Guochao Liu, Liujiang Yan, Ke Wang, Ming-Ming Cheng, and Qibin Hou. Towards stable 3d object detection. In *European conference on computer vision*. Springer, 2024.

[40] Shihao Wang, Yingfei Liu, Tiancai Wang, Ying Li, and Xiangyu Zhang. Exploring object-centric temporal modeling for efficient multi-view 3d object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3621–3631, 2023.

[41] 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. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17850–17859, 2023.

[42] Yingming Wang, Xiangyu Zhang, Tong Yang, and Jian Sun. Anchor detr: Query design for transformer-based detector. In *Proceedings of the AAAI conference on artificial intelligence*, volume 36, pages 2567–2575, 2022.

[43] Yue Wang, Vitor Campagnolo Guizilini, Tianyuan Zhang, Yilun Wang, Hang Zhao, and Justin Solomon. Detr3d: 3d object detection from multi-view images via 3d-to-2d queries. In *Conference on Robot Learning*, pages 180–191. PMLR, 2022.

[44] Zetong Yang, Li Chen, Yanan Sun, and Hongyang Li. Visual point cloud forecasting enables scalable autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14673–14684, 2024.

[45] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel Ni, and Heung-Yeung Shum. Dino: Detr with improved denoising anchor boxes for end-to-end object detection. In *The Eleventh International Conference on Learning Representations*, 2022.

[46] Yunpeng Zhang, Zheng Zhu, and Dalong Du. Occformer: Dual-path transformer for vision-based 3d semantic occupancy prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9433–9443, 2023.

[47] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. In *International Conference on Learning Representations*, 2020.

[48] Ziyue Zhu, Qiang Meng, Xiao Wang, Ke Wang, Liujiang Yan, and Jian Yang. Curricular object manipulation in lidar-based object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1125–1135, 2023.

# Appendix

## A  Broader impacts

Our work proposes an end-to-end paradigm for occupancy prediction, achieving state-of-the-art RayIoU performance with fast inference speeds. This advancement can lead to real-time and precise occupancy outcomes, which are crucial for real-world applications of autonomous driving (AD). Consequently, the most significant positive impact of our work is the enhancement of safety and response speed in AD systems.

However, the biggest negative societal impact of this work, as with any component of AD systems, is the safety concern. Autonomous driving systems are directly related to human lives, and erroneous predictions or decisions can lead to hazardous outcomes. Therefore, increasing the accuracy of occupancy outcomes and developing complementary methods to address false predictions will require substantial follow-up efforts.

## B  Licenses for involved assets

Our code is built on top of the codebase[1] provided by SparseBEV [20], which is subject to the MIT license. Our experiments are conducted on the Occ3D-nuScenes [38] which provides occupancy labels for the nuScenes dataset [2]. Occ3D-nuScenes is licensed under the MIT license, and nuScenes is licensed under the CC BY-NC-SA 4.0 license.

## C  Complexity analysis

In this part, we provide a detailed analysis of the time and space complexity involved in matching $m$ predictions with $n$ ground-truths.

**Hungarian algorithm.** The Hungarian algorithm's core involves finding augmenting paths for $\min(m, n)$ iterations. Each iteration can be visualized as an attempt to improve the current matching by finding the shortest augmenting path in the residual graph, which has a complexity of $O(\max(m, n)^2)$ using with Dijkstra's algorithm. Consequently, the time complexity for the Hungarian algorithm is $O(\min(m, n) \cdot \max(m, n)^2)$.

Meanwhile, the Hungarian Algorithm necessitates computing a cost matrix of size $m \times n$ to store the costs linked with each potential assignment. Throughout the matching process, the tracked labels and matched pairs each demand $O(\min(m, n))$ space. Hence, the final space complexity is $O(m \times n)$.

**Our method.** Our method employs the Chamfer distance loss, which involves computing pairwise distances and determining the smallest distance for each point. The first step requires a time complexity of $O(m \times n)$ and the next step requires $O(m \times n)$ as well. The assignment of semantic labels can re-use the results of previous nearest search, therefore requires no additional computations. In the end, the time complexity is $O(m \times n)$.

For each point in one set, the algorithm needs to keep track of the minimum distance to any point in the other set. This can be done using a single variable per point, resulting in $O(m)$ and $O(n)$ in the respective directions. Semantic label assignment, meanwhile, incurs a space complexity of $O(m)$. Collectively, this sums up to $O(2m + n)$.

**Comparison of the two methods.** In conclusion, when $m$ and $n$ are comparable in scale, the Hungarian algorithm exhibits time complexity of $O(n^3)$ and space complexity of $O(n^2)$, whereas our method demonstrates significantly improved efficiencies with complexities of $O(n^2)$ and $O(n)$, respectively. This represents a notable reduction in both time and space requirements, making it a more efficient solution for large-scale applications.

## D  Additional experiments.

### D.1  Comparison of Hungarian matching and our method

Tab. 6 presents a comparison of the duration and GPU utilization when matching two point clouds with the same number of points. It is evident that the Hungarian algorithm exhibits scalability issues. For instance, when the point number is 10K, it consumes approximately 24 seconds and 2,304Mb of

---

[1]https://github.com/MCG-NJU/SparseBEV

Table 6: Comparison of Hungarian algorithm and our label assignment scheme.

| Number of Points | Time (ms) | | GPU (Mb) | |
|---|---|---|---|---|
| | Hungarian Algorithm | Ours | Hungarian Algorithm | Ours |
| 100 | 0.52 | 0.12 | 39 | 14 |
| 1,000 | 78.34 | 0.13 | 81 | 14 |
| 10,000 | 24,216.35 | 1.25 | 2,304 | 15 |
| 100,000 | - | 28.85 | - | 39 |

GPU memory for a single matching. Scaling up to 100K points renders the matching infeasible due to CUDA memory constraints, even on an 80G A100 GPU.

In contrast, our label assignment method achieves remarkable efficiency, requiring only about 1.25ms and 28.85ms for 10K and 100K points, respectively. Furthermore, the GPU memory consumption during training is negligible. These findings reveal the practicality and efficacy of our label assignment approach, particularly for the occupancy prediction where point counts can easily exceed 10K.

Table 7: Configurations for different models.

| Model | $Q$ | $S$ | point number | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | s1 | s2 | s3 | s4 | s5 | s6 |
| OPUS-T | 600 | 4 | 1 | 4 | 16 | 32 | 64 | 128 |
| OPUS-S | 1200 | 2 | 1 | 4 | 8 | 16 | 32 | 64 |
| OPUS-M | 2400 | 2 | 1 | 2 | 4 | 8 | 16 | 32 |
| OPUS-L | 4800 | 2 | 1 | 2 | 4 | 8 | 16 | 16 |

## D.2 Detailed configuration for different versions.

In this section, we detail the settings of various versions of our model, as shown in Tab. 7, each tailored to prioritize different aspects of performance and speed. Our fastest model OPUS-T utilizes only 0.6K queries, with each query sampling 4 points in images. The number of predicted points are 1, 4, 16, 32, 64 and 128 for 6 stages, respectively. This configuration ensures a rapid processing time while maintaining competitive performance. Other versions of our model, such as OPUS-S, OPUS-M, OPUS-L, sample 2 points in CPS module, progressively double the number of queries and adjust the number of predicted points accordingly to balance speed and accuracy. All these models predict the same amount of points in the end.

Table 8: Performance with different points predicted.

| Model | point number | mIoU | $\text{RayIoU}_{1m}$ | $\text{RayIoU}_{2m}$ | $\text{RayIoU}_{4m}$ | RayIoU |
|---|---|---|---|---|---|---|
| OPUS-M | 64 | 28.4 | 22.2 | 29.5 | 34.8 | 28.8 |
| | 32 | 27.2 | 26.1 | 33.3 | 38.4 | 32.6 |
| | 16 | 22.8 | 28.1 | 35.3 | 40.2 | 34.5 |
| | 8 | 16.4 | 27.4 | 34.6 | 39.6 | 33.9 |

## D.3 Effects of various refined points number in last layer.

Tab. 8 assesses the impact of varying the number of predicted points in the last layer. We use OPUS-M as our model for this experiment. As shown in the table, mIoU steadily rises as the number of points increase from 8 to 64, going from 16.4 to 28.4. This trend is expected since increasing the number of points generally leads to higher mIoU by covering more voxels, as mIoU penalizes false negative (FN) heavily. However, the RayIoU results peak when model predicting 16 points and decline with further increasing points. This decline occurs partly because adding more points beyond a certain extent introduces noise, which negatively impacts RayIoU, which emphasizes first occupied voxels along the ray.

## D.4 Predictions across different distances

We report the RayIoU of FB-Occ and OPUS at different ranges in Tab. 9. It is evident that OPUS demonstrates a more pronounced advantage in nearby areas than at far distances. This could be attributed to the phenomenon pointed out by SparseOcc: dense approaches tend to overestimate the surfaces, especially in nearby areas.

Table 9: Performance across different distances.

| Model | overall | $0m \sim 20m$ | $20m \sim 40m$ | $> 40m$ |
|---|---|---|---|---|
| FB-Occ | 33.5 | 41.3 | 24.2 | 12.1 |
| OPUS-L | 41.2 | 49.10 | 31.15 | 13.73 |

# E    Additional qualitative analysis

## E.1    Differences between SparseOcc and OPUS

**View perspective of occupancy prediction.** The fundamental difference lies in the perspective of occupancy prediction. As depicted in the main draft, all previous methods, including SparseOcc [21], treat occupancy prediction as a standard classification task. OPUS, however, pioneers a set prediction viewpoint, offering a novel, elegant, and end-to-end sparsification approach.

**Multi-stage vs. end-to-end sparsification procedure.** SparseOcc generates sparse occupancy by gradually discarding voxels through multiple stages. The discarding of empty voxels at early stages is irreversible, leading to obvious cumulative errors, as illustrated in Fig. 3. Conversely, OPUS circumvents complex filtering mechanisms by directly predicting a sparse set, resulting in more coherent outcomes.

**Detailed model design.** In terms of a more detailed perspective of the structure, there are also many differences such as:

- **Query number.** In NuScene-Occ3D, SparseOcc necessitates 32K queries in its final stage. OPUS, by comparison, operates with a mere 0.6K-4.8K queries for occupancy prediction, capitalizing on its flexible nature and contributing to its fast inference pace.

- **Coarse-to-fine procedure.** SparseOcc's coarse-to-fine strategy involves progressively filtering empty voxels and subdividing occupied voxels into finer ones. In contrast, OPUS interprets coarse-to-fine as the escalation in number of predicted points across stages.

- **Learning objective.** Our learning target encompasses predicting both semantic classes and occupied locations, simultaneously. The latter is a new objective introduced by OPUS, achieved through a modified Chamfer distance loss.

## E.2    Analysis of relationships between mIoU, RayIoU and driving safety.

Our OPUS-L (8f) has achieved a state-of-the-art RayIoU of 41.17, outperforming the previous sparse model SparseOcc by 6.07 and the dense model FB-Occ by 7.7. The mIoU gap between sparse and dense methods is also reduced from 8.5 in SparseOcc to 3.0 in OPUS. However, the implications of this gap on safety remain ambiguous. This concern is particularly pertinent in the context of autonomous driving, and we would like to clarify this as follows:

**Risks of dense predictions.** The biggest issue of dense predictions is the discrepancies between evaluation metrics and real-world scenarios. As shown in Fig. 8, evaluation metrics only consider voxels within the camera mask, which is derived from camera parameters and ground truth. However, in real-world applications, we can only produce view mask based on camera intrinsics and extrinsics, failing to filtering out over-estimated voxels. From Fig. 8 and and Fig. 3, dense methods can misidentify occupied voxels, even close to the ego vehicle. These errors are overlooked during evaluation but pose significant safety hazards in real-world scenarios. In contrast, OPUS suffer much less from this issue as it does not over-estimate occupancy.

**The depth errors of OPUS is much smaller than FB-Occ.** In Fig. 9, we compare the depth errors of FB-Occ and OPUS along camera rays. OPUS demonstrates lower depth errors across all scenes, despite its relatively low mIoU performance. Given the significance of the first occupied voxel for safety, OPUS's precision in this regard enhances safety rather than detracting from it.

In conclusion, while it is necessary to minimize the mIoU gap between sparse and dense methods, our analysis indicates that mIoU might not fully represent potentially hazardous situations. Therefore, it would be more rational to take both mIoU and RayIoU into consideration for the occupancy task.
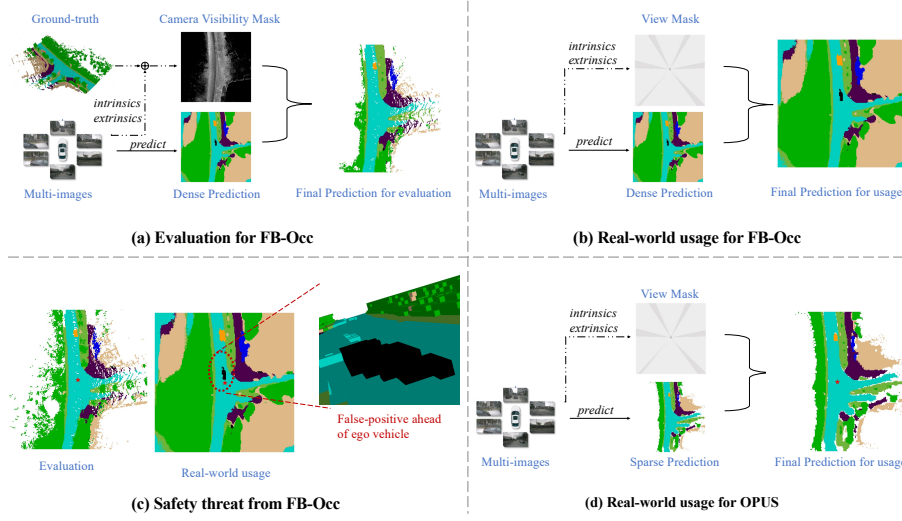
Figure 8: Illustration of safety threat due to discrepancies between evaluation metrics and real-world scenarios. (a) Before evaluation, the camera visibility mask is first generated according to camera intrinsics and extrinsics. Then, the dense prediction will be masked to get the final prediction for evaluation. (b) For real-world usage, we cannot have camera visibility reasoning without knowing the ground-truth occupancy. We can only generate the view mask from camera intrinsics and extrinsics, which fails to filter out the over-estimated voxels from dense models. (c) Plenty of false positive predictions are made close to the ego vehicle, marked by the symbol of red star. These erroneously predicted voxels are filtered during evaluating mIoU, but could cause hazardous safety issue. (d) The OPUS produces sparse occupancy predictions and suffers much less from the over-estimation. Consequently, no such safety threat occurs in this scenario. **Best viewed in color**.
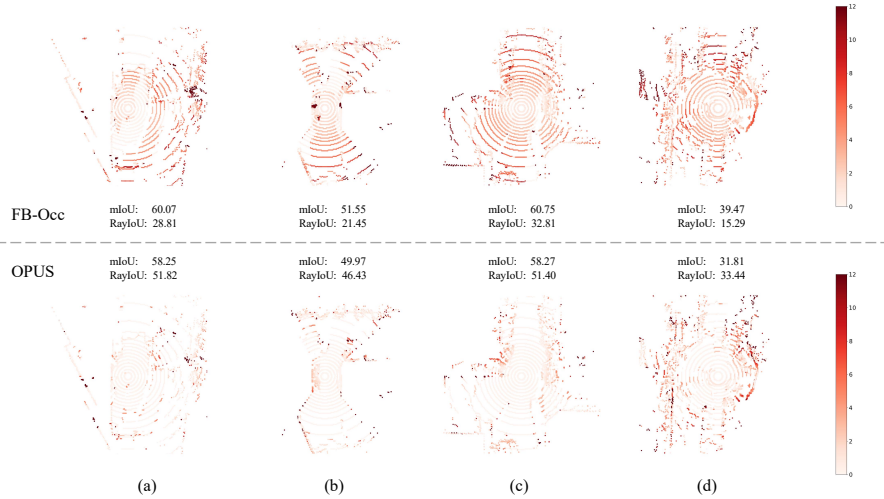


Figure 9: The predicted error maps of FB-Occ and OPUS. When compared with FB-Occ, OPUS has lower mIoU and higher RayIoU results, and achieves evidently smaller errors. **Best viewed in color**.

### E.3 Occupancy predictions of different methods.

In Fig. 10, We further provide more visualizations of occupancy predicted by FB-Occ, SparseOcc, and proposed OPUS. As shown in Fig. 3 and Fig. 10, a common OPUS failure mode is the prediction of scattered and discontinuous surfaces at long distances. Another is the presence of holes in predicted driving surface, a phenomenon also observed in SparseOcc due to the sparsity properties.
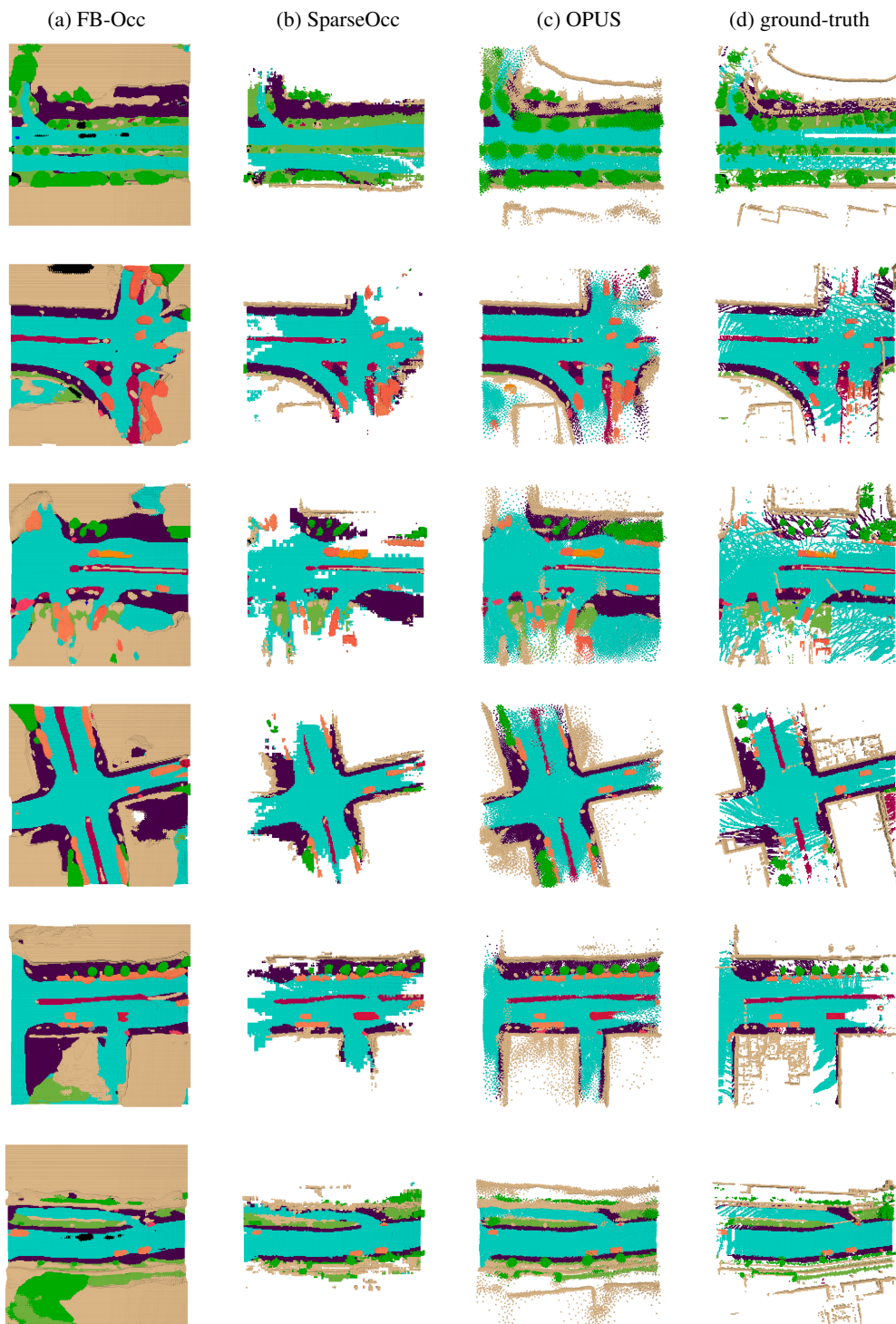


Figure 10: Visualizations of occupancy predicted by FB-Occ, SparseOcc and the proposed OPUS.