

# MARS: MAMBA-DRIVEN ADAPTIVE REORDERING SCHEME FOR SEMANTIC OCCUPANCY PREDICTION IN AUTONOMOUS DRIVING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Semantic occupancy prediction provides fundamental voxel-level scene perception for autonomous driving systems, yet the massive number of voxels poses significant computational challenges, especially for Transformer-based methods with quadratic complexity. Recently, OccMamba introduces state space models to this task, but its reliance on a handcrafted 3D-to-1D reordering scheme suffers from two critical challenges: (1) *indiscriminate processing of redundant empty voxels*, and (2) *limited adaptivity to diverse scene layouts*. To address this issue, we propose the **Mamba-driven Adaptive Reordering Scheme (MARS)** framework, replacing the static reordering scheme with an adaptive and dynamic design, facilitating modality-aware pruning of redundant empty voxels and scene-adaptive sequence of critical voxels. Specifically, we first introduce the *Adaptive Voxel Pruning (AVP)* module to tackle indiscriminate processing, which filters out redundant empty voxels and retain informative ones, thereby establishing an efficient computational foundation. Then, we present the *Dynamic Voxel Reordering (DVR)* module to address limited adaptivity, which dynamically identifies and sequences critical voxels for scene-level perception, ensuring flexible adaptivity to diverse scenarios. Extensive experiments and analyses on the OpenOccupancy dataset showcase the effectiveness and efficiency of our MARS framework, achieving superior semantic occupancy performance while reducing training memory by 19.6% and accelerating inference by 9.7%.

## 1 INTRODUCTION

Semantic occupancy prediction aims to generate a dense voxel-level semantic and geometric representation of the surrounding 3D environment Roldao et al. (2022), providing unified and structured scene comprehension foundation for autonomous driving systems. The pioneering methods, such as JS3C-Net Yan et al. (2021) and MonoScene Cao & de Charette (2022), have made progress with uni-modal partial inputs. Following researches, such as M-CONet Wang et al. (2023b) and FusionOcc Zhang et al. (2024), leverage multi-modal inputs to extract improved visual cues and contextual details, but suffer from the inherent deficiency of CNN He et al. (2016) in modeling long-range global relationships.

To overcome the inherent complexity and scale of autonomous driving scenarios, recent methods Li et al. (2023a); Tong et al. (2023); Zhang et al. (2023); Mei et al. (2024) predominantly adopt Transformer Vaswani et al. (2017) attention architectures to model long-range dependencies within the voxelized scenes. However, the massive number of voxels poses significant challenges to training efficiency due to the quadratic computational complexity of transformer-based networks. In light of this, OccMamba Li et al. (2025) first introduces Mamba Gu & Dao (2023) architecture to the semantic occupancy prediction task, processing 3D-to-1D reordered voxel sequences with linear computation complexity.

Despite its progress, the handcrafted, static 3D-to-1D reordering scheme remains problematic. We conduct a pilot study by varying the proportion of voxels fed into the Mamba blocks in OccMamba, as illustrated in Figure 1 (a). It can be observed that reducing the proportion of voxels to only 1/8 of the full size yields comparable or even slightly higher performance than processing all voxels

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

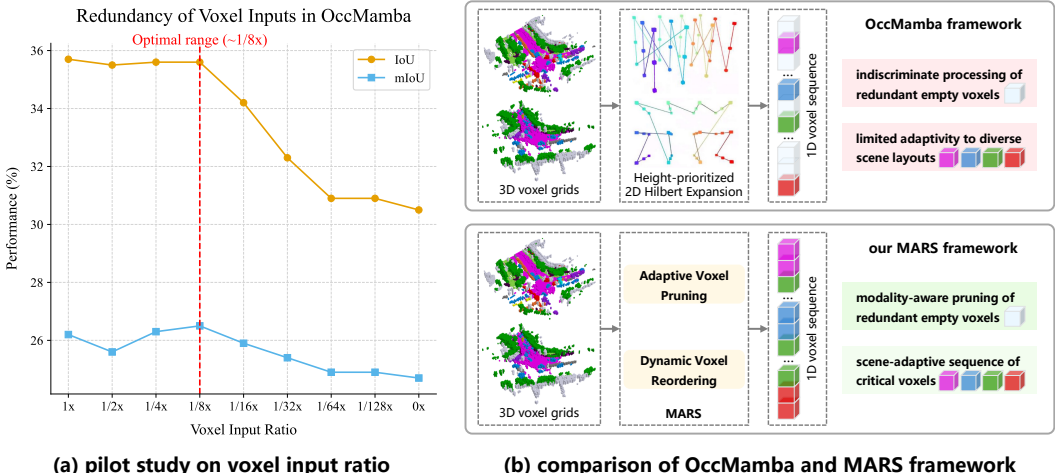


Figure 1: **Motivation for the MARS framework.** (a) Our pilot study on OccMamba shows that indiscriminately processing all voxels is suboptimal. Performance is maintained down to a 1/8 input ratio, revealing massive redundancy, yet collapses thereafter, proving the importance of a critical voxel subset. (b) This motivates our design, which replaces OccMamba’s static scheme with a two-stage adaptive process. Our MARS framework first uses Adaptive Voxel Pruning to tackle redundancy, then employs Dynamic Voxel Reordering to address the lack of adaptivity by creating a context-aware sequence of critical voxels.

throughout the scene, indicating that a large fraction of voxels contributes little to the final prediction and even impairs model performance with redundant computation. In addition, when the voxel input ratio is reduced below 1/16, the performance drops sharply, suggesting the existence of a critical subset of voxels that are essential for maintaining scene-level perception. These findings expose two major limitations of the current 3D-to-1D reordering scheme in OccMamba: (1) **Indiscriminate processing of redundant empty voxels** The handcrafted reordering scheme feeds all voxels into Mamba blocks, overlooking the fact that a large proportion (over 90%) of voxels are empty and causing redundant computation. (2) **Limited adaptivity to diverse scene layouts.** The static reordering scheme fails to adaptively prioritize the critical subset of informative voxels across diverse scene layouts, suffering from disruption of non-informative voxels.

To address these issues, we propose the **Mamba-driven Adaptive Reordering Scheme (MARS)** framework, replacing the handcrafted, static reordering scheme with an adaptive and dynamic design, as shown in Figure 1 (b). Specifically, to alleviate the redundancy of processing vast unoccupied regions, we devise the **Adaptive Voxel Pruning** module. Notice that different input modalities exhibit distinct characteristics, emphasizing different regions in the voxel space: LiDAR point clouds provide precise 3D geometric structure information but are inherently sparse, while camera images capture dense semantic cues but lack direct depth information. Therefore, we devise modality adapters to encode modality-aware meta information, followed by a lightweight 3D multi-layer perceptron network identifying informative regions and pruning redundant empty voxels. To overcome the rigidity of handcrafted reordering schemes and adapt to diverse scene layouts, we propose the **Dynamic Voxel Reordering** module. Dynamic driving environments exhibit diverse voxel space distributions, as different scenes consist of distinct arrangements of moving objects, road structures, and free space. In light of this, we design scene adapters that aggregate multi-modal features and exploit contextual voxel dependencies around informative areas, generating dynamic 3D-to-1D sequences of critical voxels with spatial continuity. This dynamic reordering scheme facilitates our MARS framework with modality-aware pruning of redundant empty voxels and scene-adaptive sequence of critical voxels, which flexibly adapts to diverse scene layouts, emphasizing informative critical voxels for state space modeling. Extensive experiments and analyses on the OpenOccupancy Wang et al. (2023b) benchmark demonstrates the effectiveness of our MARS approach.

In summary, our contributions are as follows:

- **An efficient and effective state space model architecture for semantic occupancy prediction.** We present the MARS framework, replacing the handcrafted, static 3D-to-1D reordering scheme with an adaptive and dynamic design, facilitating modality-aware pruning of redundant empty voxels and scene-adaptive sequence of critical voxels for efficient and effective semantic occupancy prediction.
- **An adaptive voxel pruning strategy.** To address the indiscriminate processing of redundant empty voxels, we leverage modality-specific adapters to filter out redundant voxels while preserving informative ones.
- **A dynamic voxel reordering scheme.** To alleviate the limited adaptivity to diverse scene layouts, we exploit contextual voxel dependencies within spatial coherent regions, dynamically prioritizing and sequencing critical voxels across diverse scenarios.

## 2 RELATED WORK

### 2.1 SEMANTIC OCCUPANCY PREDICTION

The primary objective of semantic occupancy prediction is to assess voxel-level occupancy and semantic labels of surrounding scenes. Based on different input modalities, semantic occupancy prediction methods can be broadly divided into three main streams: LiDAR-based methods, camera-based methods, and multi-modal fusion methods.

**LiDAR-based Methods.** Leveraging the natural 3D geometric structural information of LiDAR point clouds, LiDAR-based methods Rist et al. (2021); Cheng et al. (2021); Xia et al. (2023) have long been the predominant solutions to semantic occupancy prediction. UDNNet Zou et al. (2021) employs a 3D U-Net architecture to directly construct scene predictions from LiDAR point clouds. LMSCNet Roldao et al. (2020) utilizes 2D convolution networks for lightweight encoding of voxel features. SGCNet Zhang et al. (2018) splits input voxels into distinct groups, enabling sparse spatial group convolutions. JS3C-Net Yan et al. (2021) conducts knowledge fusion between semantic segmentation and occupancy prediction with point-voxel interaction. SSC-RS Mei et al. (2023) designs a multi-branch architecture to integrate semantic and geometric features hierarchically.

**Camera-based Methods.** Due to the cost-effectiveness and flexibility of camera sensors, camera-based methods Miao et al. (2023); Wang & Tong (2024); Wang et al. (2024) have garnered increasing attention with only RGB image inputs. MonoScene Cao & de Charette (2022) first samples image features along lines of sight into voxel representations. Subsequent methods employ Bird’s-Eye-View (BEV) representations Li et al. (2022b; 2023c;b) and Tri-Perspective-View (TPV) Huang et al. (2023) representations for more efficient scene modeling. Recently, transformer-based methods Zhang et al. (2023); Wei et al. (2023) adopt transformer architecture with BEV queries Li et al. (2022b), voxel queries Li et al. (2023a), instance queries Jiang et al. (2024), and context queries Yu et al. (2024) to handle the inherent complexity in autonomous driving scenes.

**Multi-modal Fusion Methods.** Towards more accurate details and robust perception, multi-modal fusion methods Gao et al. (2020); Li et al. (2022a); Wang et al. (2023a) leverage different input modalities. AICNet Li et al. (2020) integrates RGB images and depth maps with anisotropic convolution networks. Point-Painting Vora et al. (2020) utilizes image segmentation probabilities to enhance point cloud representations with rich semantics. UniSeg Liu et al. (2023a) makes use of RGB images and three views of LiDAR point cloud for panoptic semantic segmentation. BEVFusion Liu et al. (2023b) unifies multi-modal features in a shared BEV representation space to preserve semantic and geometric information.

However, most mainstream semantic occupancy prediction methods, especially transformer-based ones with quadratic computation complexity, still face significant computational challenges due to the massive number of voxels in large-scale autonomous driving scenarios.

### 2.2 STATE SPACE MODELS

State Space Models (SSMs) Gu et al. (2021); Nguyen et al. (2022) have recently attracted increasing attention as a competitive alternative to transformer Vaswani et al. (2017) architectures for model-

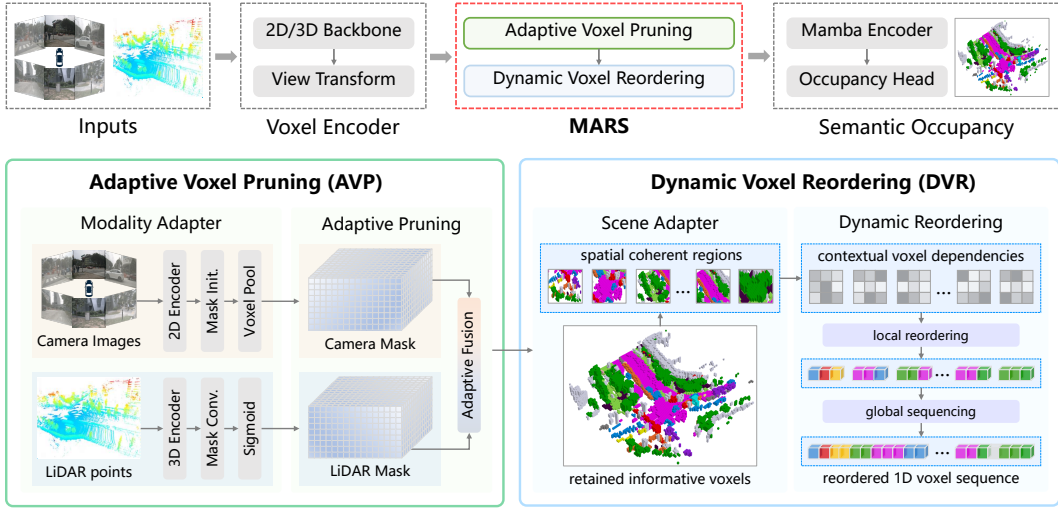


Figure 2: The overall architecture of our MARS approach. The Adaptive Voxel Pruning (AVP) module leverages modality-specific adapters to filter out redundant empty voxels while retaining informative ones. The Dynamic Voxel Reordering (DVR) module exploits contextual voxel dependencies within spatial coherent regions to dynamically prioritize and sequence critical voxels.

ing sequential data with linear computational complexity. Mamba Gu & Dao (2023), notable for processing large-scale data in linear time with selective mechanisms, further extends its variants for processing 2D and 3D data through variants like VMamba Liu et al. (2024), Vision Mamba Zhu et al. (2024), and PointMamba Liang et al. (2024). OccMamba Li et al. (2025) first designs a Mamba-based architecture for semantic occupancy prediction, where 3D voxels are transferred into 1D sequence through a height-prioritized 2D Hilbert reordering scheme. However, the handcrafted, static reordering scheme suffers from redundant processing of empty voxels and limited adaptivity to dynamic scenarios. To address this issue, we propose MARS, an adaptive and dynamic 3D-to-1D reordering scheme for improved computational efficiency and scene adaptivity.

### 3 METHODOLOGY

In this section, we first briefly review state space models (SSMs) and its application in semantic occupancy prediction (OccMamba). Then, we introduce the *Mamba-based Adaptive Reordering Scheme (MARS)* framework, an efficient state space model architecture that leverages adaptive and dynamic 3D-to-1D reordering scheme for improved computational efficiency and scene adaptivity. As illustrated in Figure 2, the overall framework of our MARS approach consists of two key components: (1) an Adaptive Voxel Pruning (AVP) module that leverages modality-specific adapters to filter out redundant voxels while preserving informative ones; (2) a Dynamic Voxel Reordering (DVR) module that exploits contextual voxel dependencies to dynamically prioritize and sequence critical voxels.

#### 3.1 PRELIMINARIES

**Problem Setup.** The primary objective of semantic occupancy prediction is to generate a dense representation of the surrounding environment, assigning both occupancy states and semantic classes to a predefined 3D voxel grids. Formally, given a set of multi-view camera images  $\mathcal{I} = \{I_i\}_{i=1}^{N_{\text{cam}}}$  and LiDAR point clouds  $\mathcal{P}$ , the task is to predict a voxel grid  $\mathcal{V} \in \mathbb{R}^{X \times Y \times Z}$ , where  $X, Y, Z$  correspond to the grid’s resolution. Each voxel is assigned a semantic label of either empty denoted by  $c_0$  or occupied by one of the predefined semantic classes in  $\mathcal{C} \in \{c_1, \dots, c_N\}$ .

**State space models.** State space models (SSMs) are inspired by the control theory and formulates sequence modeling as a structured state transition process, where the hidden state evolves according

to linear dynamical systems and is updated with input-dependent parameters. This process can be formulated as:

$$h_t = Ah_{t-1} + Bx_t, \quad y_t = Ch_t + Dx_t \quad (1)$$

where  $h_t$  is the hidden state,  $x_t$  is the input, and  $y_t$  is the output. Recent variants such as Mamba Gu & Dao (2023) enhance SSMs with selective state update and efficient parallelization, modeling long-range dependencies with linear computations complexity.

**OccMamba.** OccMamba Li et al. (2025) first introduces the Mamba-based architectures into semantic occupancy prediction, consisting of multi-modal visual encoders, an OccMamba encoder, and an occupancy head. It converts 3D voxelized grids into 1D voxel sequences through a height-prioritized 2D Hilbert reordering scheme, then employs hierarchical mamba modules and local context processors to model scene-level long-range dependencies. However, the handcrafted and static 3D-to-1D reordering scheme suffers from issues of indiscriminate processing of redundant empty voxels and limited adaptivity to diverse scene layouts, thus limiting overall performance of semantic occupancy prediction.

### 3.2 ADAPTIVE VOXEL PRUNING

To address the issue of indiscriminate processing of redundant empty voxels, we propose the Adaptive Voxel Pruning (AVP) module, leveraging modality-aware information to filter out redundant voxels while retaining informative regions. The key insight behind AVP is that considering the large proportion of empty voxels (over 90%) and complementary characteristics of different input modalities, it is crucial to leverage modality-aware metadata together with extracted features to adaptively identify informative voxels, pruning redundant ones for improved computational efficiency.

**Modality Adapter.** To fully exploit the inherent characteristics of different modalities, the modality adapter leverages modality-aware metadata to generate heuristic geometric structures within the voxel space, serving as spatial priors for redundant voxel pruning.

For camera image inputs, we leverage camera parameters and depth probabilities as metadata to compensate for the lack of explicit 3D structural information. Specifically, we first initialize the heuristic image mask  $M_{\text{heu}} \in \mathbb{R}^{1 \times H_{\text{img}} \times W_{\text{img}}}$  with all elements set to 1, indicating the potential contribution of each pixel. Then, the DepthNet CS Kumar et al. (2018) is utilized to generate pixel-wise depth probabilities  $D \in \mathbb{R}^{D \times H_{\text{img}} \times W_{\text{img}}}$ , lifting image masks along the depth dimension. Finally, the lifted masks are regularized into voxelized grid coordinates through voxel pooling Philion & Fidler (2020) with camera parameters. The above process can be formulated as follows:

$$M_{\text{cam}} = \text{VoxelPool}((M_{\text{heu}} \cdot D), K, T) \quad (2)$$

where  $K, T$  are the camera intrinsic and extrinsic matrices, and the voxel pooling operation accumulates depth-aware contributions of each pixel, generating camera voxel mask  $M_{\text{cam}} \in \mathbb{R}^{X \times Y \times Z}$ .

For LiDAR point cloud inputs, they naturally encode explicit 3D structural information with 3D coordinates and reflection intensity, which is well-suited for extracting spatial priors. To fully exploit these properties, we first regularize the input point clouds  $P$  into voxelized grids Zhou & Tuzel (2018), where each voxel aggregates both geometric coordinates and intensity values of points falling inside it. The voxelized features  $V_{\text{lidar}} \in \mathbb{R}^{C \times X \times Y \times Z}$  are then fed into a lightweight multi-layer perceptron  $\psi_{\text{lidar}}(\cdot)$  to capture geometry-aware spatial contexts:

$$M_{\text{lidar}} = \sigma(\psi_{\text{lidar}}(V_{\text{lidar}})) \quad (3)$$

where  $\sigma(\cdot)$  is the sigmoid function, and  $M_{\text{lidar}} \in \mathbb{R}^{X \times Y \times Z}$  is the LiDAR voxel mask, providing spatial priors within the voxel space based on LiDAR structural and intensity cues.

**Adaptive Pruning.** Given the modality-aware voxel masks  $M_{\text{cam}}, M_{\text{lidar}}$ , we further design an adaptive pruning module to fuse multi-modal features and filter out redundant voxels. Camera masks highlight regions aligned with projected image features, while LiDAR masks emphasize structural priors in occupied space. We generate adaptive fusion weights  $W \in \mathbb{R}^{X \times Y \times Z \times 2}$  from the concatenation of multi-modal voxel features and masks  $[V_{\text{cam}}, M_{\text{cam}}; V_{\text{lidar}}, M_{\text{lidar}}]$  using 3D convolution

with softmax on the last dimension. The fused voxel features and masks are computed as follows:

$$\begin{aligned} \mathbf{V} &= \mathbf{W}_0 \odot \mathbf{M}_{\text{cam}} \odot \mathbf{V}_{\text{cam}} + \mathbf{W}_1 \odot \mathbf{M}_{\text{lidar}} \odot \mathbf{V}_{\text{lidar}} \\ \mathbf{M} &= \mathbf{W}_0 \odot \mathbf{M}_{\text{cam}} + \mathbf{W}_1 \odot \mathbf{M}_{\text{lidar}} \end{aligned} \quad (4)$$

where  $\odot$  is the element-wise multiplication,  $\mathbf{V}$  is the fused voxel features, and  $\mathbf{M}$  is the adaptive voxel mask with voxel-wise occupancy scores based on multi-modal spatial priors, facilitating redundant voxel pruning with confidence threshold  $\mathbf{1}_{M>\theta}$ .

### 3.3 DYNAMIC VOXEL REORDERING

To address the issue of limited adaptivity to diverse scene layouts, we propose the Dynamic Voxel Reordering (DVR) module, exploiting contextual voxel dependencies to dynamically prioritize and sequence critical voxels across diverse scenarios. The design rationale of DVR is that since diverse driving environments consist of distinct voxel space distributions, it is essential to exploit scene-aware voxel dependencies and dynamically reorder critical voxels, enhancing the adaptivity to diverse scenarios. DVR takes the fused features  $\mathbf{V}$  and pruning mask  $\mathbf{M}$  from AVP as input.

**Scene Adapter.** To generate a reordering scheme that is aware of the overall scene layout, the scene adapter is designed to discover and summarize spatially coherent regions of interest. We first employ a connected components labeling algorithm on the pruning mask to partition the set of informative voxels  $\mathbf{1}_{M>\theta}$  into  $K$  disjoint salient regions  $\{\mathbf{M}_1, \mathbf{M}_2, \dots, \mathbf{M}_K\}$ . Then for each identified region  $\mathbf{M}_k$ , we generate a representative context vector by applying an aggregation function over its constituent voxels. Specifically, we use average pooling for global aggregation:

$$\mathbf{m}_k = \text{AvgPool}(\mathbf{V}[\mathbf{M}_k]) \quad (5)$$

where  $\mathbf{m}_k$  is the representative vector of informative voxels within region  $\mathbf{M}_k$ . This process yields a set of regional context vectors, each summarizing a semantically and spatially cohesive part of the scene, thereby providing richer and more structured spatial priors for dynamic reordering.

**Dynamic Reordering.** With the partition of spatially coherent regions and the extraction of regional representative vectors, the dynamic reordering process generates a scene-adaptive 1D voxel sequence, where the importance score of each voxel is conditioned on the coherent regions it belongs to. Specifically, for each voxel  $v_i$  located in region  $\mathbf{M}_k$ , we concatenate its feature with the corresponding regional representative vector. The combined representations are then fed into a shared MLP,  $\phi_{\text{order}}$ , predicting voxel-wise reordering score  $s_i$ :

$$s_i = \phi_{\text{order}}(\text{Concat}(v_i, \mathbf{m}_k)), \quad \text{where } i \in \mathbf{M}_k \quad (6)$$

Voxels within each spatial coherent region are locally reordered based on their reordering scores. Then, the reordered sequences are globally concatenated together and fed into the Mamba encoder blocks. By conditioning the importance score on regional context, this scheme encourages the model to group and prioritize related information. For example, all voxels belonging to a "pedestrian" cluster might be ranked closely together. This enhances the locality of the input sequence and allows the state space model to build a more coherent understanding of distinct scene elements, leading to superior adaptivity across diverse scenarios.

### 3.4 TRAINING OBJECTIVE.

Following Li et al. (2025), we adopt the cross-entropy loss  $\mathcal{L}_{ce}$  for classification, the lovasz-softmax loss  $\mathcal{L}_{\text{lovasz}}$  Berman et al. (2018) for semantic segmentation, the scene-class affinity loss  $\mathcal{L}_{\text{scal}}$  Cao & de Charette (2022) for spatial alignments, and the depth supervision loss  $\mathcal{L}_{\text{depth}}$  Li et al. (2023c) for depth estimation. The final training objective is formulated as follows:

$$\mathcal{L} = \mathcal{L}_{ce} + \mathcal{L}_{\text{lovasz}} + \mathcal{L}_{\text{scal}} + \mathcal{L}_{\text{depth}} \quad (7)$$

## 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETUP

**Datasets.** We evaluate our MARS approach on the OpenOccupancy Wang et al. (2023b) dataset. OpenOccupancy extends the popular nuScenes Caesar et al. (2020) dataset with dense semantic

Table 1: Quantitative comparisons on the OpenOccupancy Wang et al. (2023b) validation set with v0.0 annotations. C, D, L denote camera, depth and LiDAR, respectively. Best results are highlighted in **bold**, and second-best results are underlined.

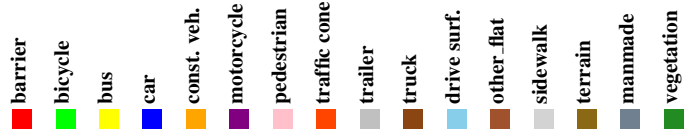
Method	Input Modality	IoU																		
		IoU	mIoU	barrier	bicycle	bus	car	const. veh.	motorcycle	pedestrian	traffic cone	trailer	truck	drive surf.	other flat	sidewalk	terrain	manmade	vegetation	
MonoScene 2022	C	18.4	6.9	7.1	3.9	9.3	7.2	5.6	3.0	5.9	4.4	4.9	4.2	14.9	6.3	7.9	7.4	10.0	7.6	
TPVFormer 2023	C	15.3	7.8	9.3	4.1	11.3	10.1	5.2	4.3	5.9	5.3	6.8	6.5	13.6	9.0	8.3	8.0	9.2	8.2	
SparseOcc 2024	C	21.8	14.1	16.1	9.3	15.1	18.6	7.3	9.4	11.2	9.4	7.2	13.0	31.8	21.7	20.7	18.8	6.1	10.6	
3DSketch 2020	C&D	25.6	10.7	12.0	5.1	10.7	12.4	6.5	4.0	5.0	6.3	8.0	7.2	21.8	14.8	13.0	11.8	12.0	21.2	
AICNet 2020	C&D	23.8	10.6	11.5	4.0	11.8	12.3	5.1	3.8	6.2	6.0	8.2	7.5	24.1	13.0	12.8	11.5	11.6	20.2	
LMSCNet 2020	L	27.3	11.5	12.4	4.2	12.8	12.1	6.2	4.7	6.2	6.3	8.8	7.2	24.2	12.3	16.6	14.1	13.9	22.2	
JS3C-Net 2021	L	30.2	12.5	14.2	3.4	13.6	12.0	7.2	4.3	7.3	6.8	9.2	9.1	27.9	15.3	14.9	16.2	14.0	24.9	
M-CONet 2023b	C&L	29.5	20.1	23.3	13.3	21.2	24.3	15.3	15.9	18.0	13.3	15.3	20.7	33.2	21.0	22.5	21.5	19.6	23.2	
Co-Occ 2024	C&L	30.6	21.9	26.5	16.8	22.3	27.0	10.1	20.9	20.7	14.5	16.4	21.6	36.9	23.5	5.5	23.7	20.5	23.5	
OccMamba 2025	C&L	35.7	26.2	30.2	20.5	<b>26.5</b>	29.5	<b>18.8</b>	26.0	23.7	19.9	<b>20.6</b>	<u>25.4</u>	<b>38.4</b>	<u>26.5</u>	<u>27.0</u>	<u>26.6</u>	<u>28.9</u>	30.5	
<b>MARS (ours)</b>	C&L	<b>36.2</b>	<b>27.1</b>	<b>31.8</b>	<b>23.0</b>	<u>26.3</u>	<b>31.0</b>	<u>17.4</u>	<b>27.8</b>	<b>27.8</b>	<b>20.1</b>	<u>20.4</u>	<b>26.1</b>	<b>39.5</b>	<b>27.2</b>	<b>27.7</b>	<b>27.1</b>	<b>29.0</b>	<b>30.9</b>	

Table 2: Efficiency evaluation results against SSM-based method OccMamba Li et al. (2025) and Transformer-based method M-CONet Wang et al. (2023b) on the OpenOccupancy validation set, taking both camera images and LiDAR point clouds as inputs.

Method	Training Memory (GB)	Inference Time (ms)	Input Voxel	mIoU (%)
M-CONet	37.3	2,231	-	20.1
OccMamba-384	37.7	2,401	163,840	26.2
OccMamba-128	23.1	2,027	163,840	25.2
MARS-384	30.3	2,168	~20,480	27.1
MARS-128	22.3	1,895	~20,480	25.7

occupancy annotations, which comprises 700 training sequences and 150 validation sequences, with a total annotation of 17 semantic classes. The semantic occupancy labels are represented within  $512 \times 512 \times 40$  voxelized grids, with a voxel resolution of  $0.2m$ .

**Evaluation Metrics.** Following Li et al. (2025), we employ the Intersection over Union (IoU) of occupied voxels as the evaluation metric for the task of class-agnostic scene scene (SC). Additionally, we report the mean Intersection over Union (mIoU) across all semantic categories to measure the performance of the semantic scene completion (SSC):

$$IoU = \frac{TP}{TP + FP + FN}, \quad mIoU = \frac{1}{C} \sum_{c=1}^C \frac{TP_c}{TP_c + FP_c + FN_c} \quad (8)$$

where  $TP, FP, FN$  represent the number of true positive, false positive, and false negative occupancy predictions, and  $C$  stands for the total number of classes.

**Implementation Details.** We employ the ResNet-50 He et al. (2016) network as the image backbone. The Mamba Encoder and Decoder maintains the architecture in OccMamba Li et al. (2025), which contains four groups and each group consists of two Mamba blocks. For the OpenOccupancy dataset, we follow the official input settings Wang et al. (2023b), where six surround view images are used as camera inputs, together with a fusion of ten frames of LiDAR points covering the range of  $[-51.2m \sim 51.2m, -51.2m \sim 51.2m, -2.0m \sim 6.0m]$ . The feature dimension within the Mamba blocks is set to 384. The threshold  $\theta$  for adaptive pruning is set default to 0.5. We train MARS for 20 epochs on 8 NVIDIA A6000 GPUs, with a total batch size of 8. The AdamW Loshchilov & Hutter (2017) optimizer is adopted with an initial learning rate of  $5e-4$  and a weight decay of  $1e-2$ .

Table 3: Ablation study on the OpenOccupancy validation set, investigating the effectiveness of different architectural components of our MARS approach with different input modalities. C denotes camera inputs and L denotes LiDAR inputs, respectively.

Variants		C		L		C&L	
AVP	DVR	IoU	mIoU	IoU	mIoU	IoU	mIoU
		19.3	12.1	35.4	22.7	35.7	26.2
✓		19.8	12.6	35.7	23.0	35.9	26.3
✓	✓	20.4	13.0	36.5	23.6	36.2	27.1

## 4.2 MAIN RESULTS

As shown in Table 1, we compare MARS with existing state-of-the-art methods on the OpenOccupancy validation set. It can be observed that our proposed MARS approach achieves a superior performance of **36.2%** IoU and **27.1%** mIoU, demonstrating its effectiveness over the strong OccMamba baseline. It is crucial to highlight that MARS achieves this result while operating on a dynamically pruned subset of voxels for improved computational efficiency. Table 2 presents the efficiency evaluations results of our MARS against SSM-based method OccMamba and Transformer-based method M-CONet, where the networks are trained on 8 NVIDIA A6000 GPUs and perform inference on a single NVIDIA A6000 GPU. Compared to OccMamba’s processing all  $128 \times 128 \times 10 = 163,840$  voxels indiscriminately, our MARS feeds only  $\sim$ **12.5%** of the total voxels into the Mamba blocks, reducing the training memory by **19.6%** (from 37.7GB to 30.3GB) and accelerating inference time by **9.7%** (from 2,401ms to 2,168ms). The above concurrent improvements in both effectiveness and efficiency are the direct outcome of our targeted architectural design, which systematically addresses the two key issues identified in our pilot study. Specifically, the Adaptive Voxel Pruning (AVP) module directly confronts the problem of indiscriminate processing. By intelligently filtering out vast, non-informative regions, it establishes a computationally efficient foundation, ensuring that the model’s resources are not wasted on empty space. Building upon this sparse yet salient representation, the Dynamic Voxel Reordering (DVR) module then tackles the challenge of limited adaptivity. It analyzes the contextual relationships between the retained voxels and sequences them into a coherent, scene-aware stream that is optimal for state space modeling. This synergy—where AVP first provides the voxel efficiency and DVR then unlocks the performance from that critical information—validates our adaptive approach.

## 4.3 ABLATION STUDY

To further investigate the effectiveness of our MARS approach and different components, we conduct ablation experiments on the OpenOccupancy validation set as follows:

**Ablation on Network Components.** As shown in Table 3, we take OccMamba Li et al. (2025) as the baseline method and present the results of ablation experiments with different modality inputs. We take OccMamba Li et al. (2025) as our baseline, where a handcrafted, static height-prioritized 2D Hilbert expansion is utilized to process all voxels with Mamba blocks. By integrating the Adaptive Voxel Pruning (AVP) module, we effectively filter out a majority of the redundant, empty voxels while retaining the informative ones. AVP provides the voxel foundation for further efficiency and effectiveness gains, and still achieves competitive performance with different modality inputs. This underscores AVP’s ability to adaptively preserve informative voxels for scene understanding. Subsequently, we introduce the Dynamic Voxel Reordering (DVR) module on top of AVP to dynamically prioritize and sequence critical voxels by exploits contextual voxel dependencies, further improves the semantic occupancy prediction performance with 1.1% IoU and 0.9% mIoU improvements for camera inputs, 1.1% IoU and 0.9% mIoU improvements for LiDAR inputs, and 0.5% IoU and 0.9% mIoU improvements for multi-modal inputs, respectively. These improvements validate that the synergy of first pruning redundancy (AVP) and then intelligently structuring the remaining information (DVR) is key to our framework’s success.

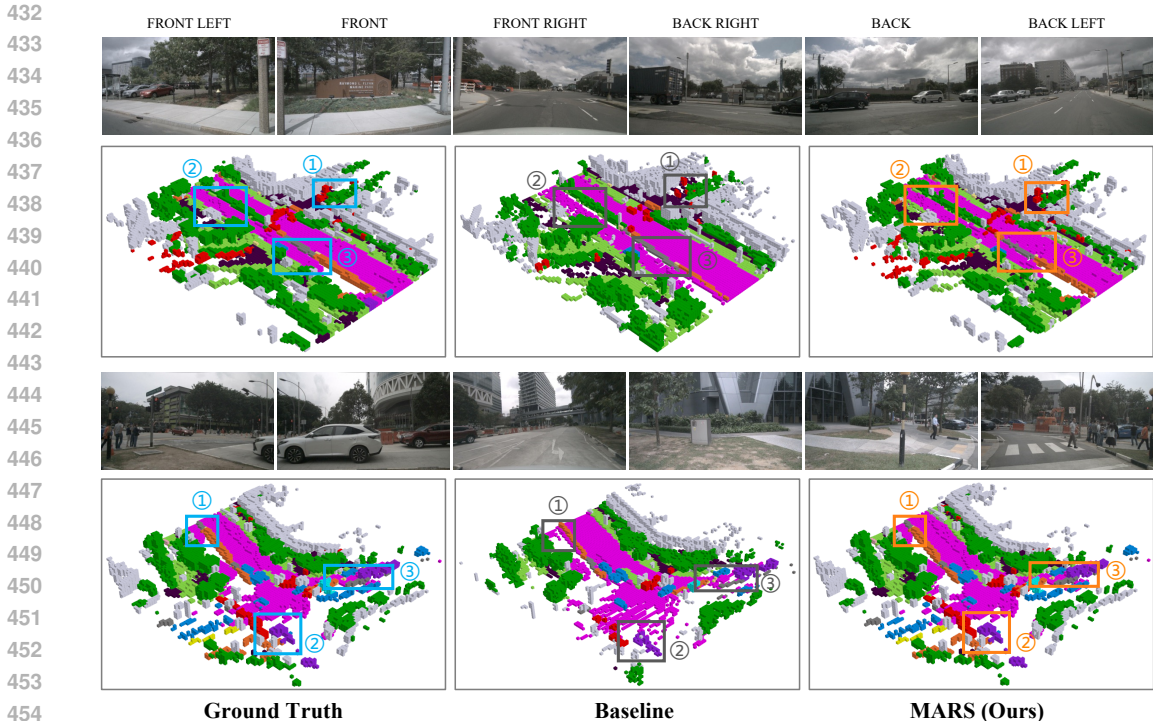


Figure 3: Visualization results on the OpenOccupancy validation set. The occupancy ground truth is outlined with blue boxes. While black boxes indicate erroneous occupancy predictions of the baseline method, and orange boxes highlight more accurate predictions by our MARS approach. Better viewed when zoomed in.

#### 4.4 VISUALIZATIONS

Figure 3 demonstrates the visualization results from the OpenOccupancy validation set. The surround-view input images are illustrated in the first and third lines. In the first row, the occupancy ground truth is outlined with blue boxes. The second row presents the occupancy predictions generated by the baseline method, where false predictions are indicated with black boxes. While the third row displays the results of our MARS approach, and orange boxes highlight our refinement for more accurate occupancy predictions. It can be observed that the handcrafted, static reordering scheme of OccMamba struggles to maintain the spatial continuity of large entities, while indiscriminate processing of redundant voxels introduces confusing semantics. In contrast, our MARS generates more complete and accurate semantic occupancy predictions, demonstrating the effectiveness of AVP’s adaptively filtering out redundant empty voxels and DVR’s dynamically sequencing critical voxels for improved semantic occupancy prediction performance.

### 5 CONCLUSION

In this work, we addressed the critical limitations inherent in the static reordering schemes used by state space models for semantic occupancy prediction: the indiscriminate processing of redundant voxels and limited adaptivity to diverse scene layouts. To address these issues, we introduce MARS, a Mamba-driven Adaptive Reordering Scheme, replacing this static design with an adaptive and dynamic approach. MARS is composed of two synergistic modules: the Adaptive Voxel Pruning (AVP) module leverages multi-modal priors to adaptively prune redundant empty voxels, and the Dynamic Voxel Reordering (DVR) module analyzes the remaining informative voxels to generate a scene-aware, contextually prioritized 1D voxel sequence. Extensive experiments on the OpenOccupancy benchmark demonstrate that MARS achieves superior semantic occupancy prediction performance, while delivering substantial improvements in computational efficiency, significantly reducing both training memory and inference latency compared to the baseline.

## REFERENCES

- 486  
487  
488 Maxim Berman, Amal Rannen Triki, and Matthew B Blaschko. The lovász-softmax loss: A tractable  
489 surrogate for the optimization of the intersection-over-union measure in neural networks. In *Pro-*  
490 *ceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4413–4421,  
491 2018.
- 492 Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush  
493 Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for  
494 autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*  
495 *recognition*, pp. 11621–11631, 2020.
- 496  
497 Anh-Quan Cao and Raoul de Charette. Monoscene: Monocular 3d semantic scene completion.  
498 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
499 3991–4001, 2022.
- 500 Xiaokang Chen, Kwan-Yee Lin, Chen Qian, Gang Zeng, and Hongsheng Li. 3d sketch-aware se-  
501 mantic scene completion via semi-supervised structure prior. In *Proceedings of the IEEE/CVF*  
502 *Conference on Computer Vision and Pattern Recognition*, pp. 4193–4202, 2020.
- 503  
504 Ran Cheng, Christopher Agia, Yuan Ren, Xinhai Li, and Liu Bingbing. S3cnet: A sparse semantic  
505 scene completion network for lidar point clouds. In *Conference on Robot Learning*, pp. 2148–  
506 2161. PMLR, 2021.
- 507  
508 Arun CS Kumar, Suchendra M Bhandarkar, and Mukta Prasad. Depthnet: A recurrent neural net-  
509 work architecture for monocular depth prediction. In *Proceedings of the IEEE Conference on*  
510 *Computer Vision and Pattern Recognition Workshops*, pp. 283–291, 2018.
- 511  
512 Jing Gao, Peng Li, Zhikui Chen, and Jianing Zhang. A survey on deep learning for multimodal data  
513 fusion. *Neural computation*, 32(5):829–864, 2020.
- 514  
515 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*  
516 *preprint arXiv:2312.00752*, 2023.
- 517  
518 Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured  
519 state spaces. *arXiv preprint arXiv:2111.00396*, 2021.
- 520  
521 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
522 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.  
523 770–778, 2016.
- 524  
525 Yuanhui Huang, Wenzhao Zheng, Yunpeng Zhang, Jie Zhou, and Jiwen Lu. Tri-perspective view  
526 for vision-based 3d semantic occupancy prediction. In *Proceedings of the IEEE/CVF Conference*  
527 *on Computer Vision and Pattern Recognition*, pp. 9223–9232, 2023.
- 528  
529 Haoyi Jiang, Tianheng Cheng, Naiyu Gao, Haoyang Zhang, Tianwei Lin, Wenyu Liu, and Xing-  
530 gang Wang. Symphonize 3d semantic scene completion with contextual instance queries. In *Pro-*  
531 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 20258–  
532 20267, 2024.
- 533  
534 Heng Li, Yuenan Hou, Xiaohan Xing, Yuexin Ma, Xiao Sun, and Yanyong Zhang. Occmamba:  
535 Semantic occupancy prediction with state space models. In *Proceedings of the Computer Vision*  
536 *and Pattern Recognition Conference*, pp. 11949–11959, 2025.
- 537  
538 Jie Li, Kai Han, Peng Wang, Yu Liu, and Xia Yuan. Anisotropic convolutional networks for 3d  
539 semantic scene completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
*and Pattern Recognition*, pp. 3351–3359, 2020.
- 540  
541 Xin Li, Botian Shi, Yuenan Hou, Xingjiao Wu, Tianlong Ma, Yikang Li, and Liang He. Homoge-  
542 neous multi-modal feature fusion and interaction for 3d object detection. In *European Conference*  
543 *on Computer Vision*, pp. 691–707. Springer, 2022a.

- 540 Yiming Li, Zhiding Yu, Christopher Choy, Chaowei Xiao, Jose M Alvarez, Sanja Fidler, Chen Feng,  
541 and Anima Anandkumar. Voxformer: Sparse voxel transformer for camera-based 3d semantic  
542 scene completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern  
543 Recognition*, pp. 9087–9098, 2023a.
- 544 Yinhao Li, Han Bao, Zheng Ge, Jinrong Yang, Jianjian Sun, and Zeming Li. Bevstereo: Enhancing  
545 depth estimation in multi-view 3d object detection with temporal stereo. In *Proceedings of the  
546 AAAI Conference on Artificial Intelligence*, volume 37, pp. 1486–1494, 2023b.
- 548 Yinhao Li, Zheng Ge, Guanyi Yu, Jinrong Yang, Zengran Wang, Yukang Shi, Jianjian Sun, and Zem-  
549 ing Li. Bevdepth: Acquisition of reliable depth for multi-view 3d object detection. In *Proceedings  
550 of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 1477–1485, 2023c.
- 551 Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng Dai.  
552 Bevformer: Learning bird’s-eye-view representation from multi-camera images via spatiotemporal  
553 transformers. In *European conference on computer vision*, pp. 1–18. Springer, 2022b.
- 555 Dingkang Liang, Xin Zhou, Wei Xu, Xingkui Zhu, Zhikang Zou, Xiaoqing Ye, Xiao Tan, and Xiang  
556 Bai. Pointmamba: A simple state space model for point cloud analysis. *Advances in neural  
557 information processing systems*, 37:32653–32677, 2024.
- 558 Youquan Liu, Runnan Chen, Xin Li, Lingdong Kong, Yuchen Yang, Zhaoyang Xia, Yeqi Bai, Xinge  
559 Zhu, Yuexin Ma, Yikang Li, et al. Uniseg: A unified multi-modal lidar segmentation network and  
560 the openpcseg codebase. In *Proceedings of the IEEE/CVF International Conference on Computer  
561 Vision*, pp. 21662–21673, 2023a.
- 563 Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, Jianbin  
564 Jiao, and Yunfan Liu. Vmamba: Visual state space model. *Advances in neural information  
565 processing systems*, 37:103031–103063, 2024.
- 566 Zhijian Liu, Haotian Tang, Alexander Amini, Xinyu Yang, Huizi Mao, Daniela L Rus, and Song  
567 Han. Bevfusion: Multi-task multi-sensor fusion with unified bird’s-eye view representation. In  
568 *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 2774–2781. IEEE,  
569 2023b.
- 571 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint  
572 arXiv:1711.05101*, 2017.
- 573 Jianbiao Mei, Yu Yang, Mengmeng Wang, Tianxin Huang, Xuemeng Yang, and Yong Liu. Ssc-rs:  
574 Elevate lidar semantic scene completion with representation separation and bev fusion. In *2023  
575 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 1–8. IEEE,  
576 2023.
- 578 Jianbiao Mei, Yu Yang, Mengmeng Wang, Junyu Zhu, Jongwon Ra, Yukai Ma, Laijian Li, and  
579 Yong Liu. Camera-based 3d semantic scene completion with sparse guidance network. *IEEE  
580 Transactions on Image Processing*, 2024.
- 581 Ruihang Miao, Weizhou Liu, Mingrui Chen, Zheng Gong, Weixin Xu, Chen Hu, and Shuchang  
582 Zhou. Occdepth: A depth-aware method for 3d semantic scene completion. *arXiv preprint  
583 arXiv:2302.13540*, 2023.
- 585 Eric Nguyen, Karan Goel, Albert Gu, Gordon Downs, Preey Shah, Tri Dao, Stephen Baccus, and  
586 Christopher Ré. S4nd: Modeling images and videos as multidimensional signals with state spaces.  
587 *Advances in neural information processing systems*, 35:2846–2861, 2022.
- 588 Jingyi Pan, Zipeng Wang, and Lin Wang. Co-occ: Coupling explicit feature fusion with volume  
589 rendering regularization for multi-modal 3d semantic occupancy prediction. *IEEE Robotics and  
590 Automation Letters*, 2024.
- 592 Jonah Philion and Sanja Fidler. Lift, splat, shoot: Encoding images from arbitrary camera rigs  
593 by implicitly unprojecting to 3d. In *Computer Vision–ECCV 2020: 16th European Conference,  
Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16*, pp. 194–210. Springer, 2020.

- 594 Christoph B Rist, David Emmerichs, MarkusENZweiler, and Dariu M Gavrilă. Semantic scene com-  
595 pletion using local deep implicit functions on lidar data. *IEEE transactions on pattern analysis*  
596 *and machine intelligence*, 44(10):7205–7218, 2021.
- 597 Luis Roldao, Raoul de Charette, and Anne Verroust-Blondet. Lmscnet: Lightweight multiscale 3d  
598 semantic completion. In *2020 International Conference on 3D Vision (3DV)*, pp. 111–119. IEEE,  
599 2020.
- 600 Luis Roldao, Raoul De Charette, and Anne Verroust-Blondet. 3d semantic scene completion: A  
601 survey. *International Journal of Computer Vision*, 130(8):1978–2005, 2022.
- 602 Pin Tang, Zhongdao Wang, Guoqing Wang, Jilai Zheng, Xiangxuan Ren, Bailan Feng, and Chao Ma.  
603 Sparseocc: Rethinking sparse latent representation for vision-based semantic occupancy predic-  
604 tion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,  
605 pp. 15035–15044, 2024.
- 606 Wenwen Tong, Chonghao Sima, Tai Wang, Li Chen, Silei Wu, Hanming Deng, Yi Gu, Lewei Lu,  
607 Ping Luo, Dahua Lin, et al. Scene as occupancy. In *Proceedings of the IEEE/CVF International*  
608 *Conference on Computer Vision*, pp. 8406–8415, 2023.
- 609 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
610 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*  
611 *tion processing systems*, 30, 2017.
- 612 Sourabh Vora, Alex H Lang, Bassam Helou, and Oscar Beijbom. Pointpainting: Sequential fusion  
613 for 3d object detection. In *Proceedings of the IEEE/CVF conference on computer vision and*  
614 *pattern recognition*, pp. 4604–4612, 2020.
- 615 Haiyang Wang, Hao Tang, Shaoshuai Shi, Aoxue Li, Zhenguo Li, Bernt Schiele, and Liwei Wang.  
616 Unitr: A unified and efficient multi-modal transformer for bird’s-eye-view representation. In *Pro-*  
617 *ceedings of the IEEE/CVF international conference on computer vision*, pp. 6792–6802, 2023a.
- 618 Song Wang, Jiawei Yu, Wentong Li, Wenyu Liu, Xiaolu Liu, Junbo Chen, and Jianke Zhu. Not  
619 all voxels are equal: Hardness-aware semantic scene completion with self-distillation. In *Pro-*  
620 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14792–  
621 14801, 2024.
- 622 Xiaofeng Wang, Zheng Zhu, Wenbo Xu, Yunpeng Zhang, Yi Wei, Xu Chi, Yun Ye, Dalong Du, Ji-  
623 wen Lu, and Xingang Wang. Openoccupancy: A large scale benchmark for surrounding semantic  
624 occupancy perception. In *Proceedings of the IEEE/CVF International Conference on Computer*  
625 *Vision*, pp. 17850–17859, 2023b.
- 626 Yu Wang and Chao Tong. H2gformer: Horizontal-to-global voxel transformer for 3d semantic  
627 scene completion. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38,  
628 pp. 5722–5730, 2024.
- 629 Yi Wei, Linqing Zhao, Wenzhao Zheng, Zheng Zhu, Jie Zhou, and Jiwen Lu. Surroundocc: Multi-  
630 camera 3d occupancy prediction for autonomous driving. In *Proceedings of the IEEE/CVF Inter-*  
631 *national Conference on Computer Vision*, pp. 21729–21740, 2023.
- 632 Zhaoyang Xia, Youquan Liu, Xin Li, Xinge Zhu, Yuexin Ma, Yikang Li, Yuenan Hou, and Yu Qiao.  
633 Scpnet: Semantic scene completion on point cloud. In *Proceedings of the IEEE/CVF conference*  
634 *on computer vision and pattern recognition*, pp. 17642–17651, 2023.
- 635 Xu Yan, Jiantao Gao, Jie Li, Ruimao Zhang, Zhen Li, Rui Huang, and Shuguang Cui. Sparse single  
636 sweep lidar point cloud segmentation via learning contextual shape priors from scene completion.  
637 In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 3101–3109,  
638 2021.
- 639 Zhu Yu, Runmin Zhang, Jiacheng Ying, Junchen Yu, Xiaohai Hu, Lun Luo, Si-Yuan Cao, and Hui-  
640 Liang Shen. Context and geometry aware voxel transformer for semantic scene completion. *arXiv*  
641 *preprint arXiv:2405.13675*, 2024.

648 Jiahui Zhang, Hao Zhao, Anbang Yao, Yurong Chen, Li Zhang, and Hongen Liao. Efficient seman-  
649 tic scene completion network with spatial group convolution. In *Proceedings of the European*  
650 *Conference on Computer Vision (ECCV)*, pp. 733–749, 2018.  
651

652 Shuo Zhang, Yupeng Zhai, Jilin Mei, and Yu Hu. Fusionocc: Multi-modal fusion for 3d occupancy  
653 prediction. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pp. 787–  
654 796, 2024.

655 Yunpeng Zhang, Zheng Zhu, and Dalong Du. Occformer: Dual-path transformer for vision-based  
656 3d semantic occupancy prediction. In *Proceedings of the IEEE/CVF International Conference on*  
657 *Computer Vision*, pp. 9433–9443, 2023.  
658

659 Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection.  
660 In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4490–  
661 4499, 2018.

662 Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang. Vi-  
663 sion mamba: Efficient visual representation learning with bidirectional state space model. *arXiv*  
664 *preprint arXiv:2401.09417*, 2024.

665 Hao Zou, Xuemeng Yang, Tianxin Huang, Chujuan Zhang, Yong Liu, Wanlong Li, Feng Wen, and  
666 Hongbo Zhang. Up-to-down network: Fusing multi-scale context for 3d semantic scene comple-  
667 tion. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp.  
668 16–23. IEEE, 2021.  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701