

# Collaborative Learning of Local 3D Occupancy Prediction and Versatile Global Occupancy Mapping

Shanshuai Yuan<sup>1,2</sup>, Julong Wei<sup>1</sup>, Muer Tie<sup>1</sup>, Xiangyun Ren<sup>2</sup>, Zhongxue Gan<sup>1</sup>, and Wenchao Ding<sup>1\*</sup>

**Abstract**—Vision-based 3D semantic occupancy prediction is vital for autonomous driving, enabling unified modeling of static infrastructure and dynamic agents. Global occupancy maps serve as long-term memory priors, providing valuable historical context that enhances local perception. This is particularly important in challenging scenarios such as occlusion or poor illumination, where current and nearby observations may be unreliable or incomplete. Priors aggregated from previous traversals under better conditions help fill gaps and enhance the robustness of local 3D occupancy prediction. In this paper, we propose Long-term Memory Prior Occupancy (LMPOcc), a plug-and-play framework that incorporates global occupancy priors to boost local prediction and simultaneously updates global maps with new observations. To realize the information gain from global priors, we design an efficient and lightweight Current-Prior Fusion module that adaptively integrates prior and current features. Meanwhile, we introduce a model-agnostic prior format to enable continual updating of global occupancy and ensure compatibility across diverse prediction baselines. LMPOcc achieves state-of-the-art local occupancy prediction performance validated on the Occ3D-nuScenes benchmark, especially on static semantic categories. Furthermore, we verify LMPOcc’s capability to build large-scale global occupancy maps through multi-vehicle crowdsourcing, and utilize occupancy-derived dense depth to support the construction of 3D open-vocabulary maps. Our method opens up a new paradigm for continuous global information updating and storage, paving the way towards more comprehensive and scalable scene understanding in large outdoor environments.

## I. INTRODUCTION

Vision-based 3D semantic occupancy prediction is fundamental for autonomous driving systems, enabling precise and unified understanding of both static infrastructure and dynamic agents. However, perception quality often varies significantly in complex real-world environments due to dynamic factors such as weather, illumination changes, and occlusions. These factors cause local sensor observations to be unreliable or incomplete, limiting robust 3D occupancy prediction. Existing works address this issue by fusing temporal information through techniques like BEV feature

This work was supported in part by the National Natural Science Foundation of China (NSFC) under Grant 62403142, in part by the Science and Technology Commission of Shanghai Municipality under Grant 24511103100, and in part by the State Key Laboratory of Intelligent Vehicle Safety Technology under Grant IVSTSKL-202317. The computations in this research were performed using the CFFF platform of Fudan University.

Project page: <https://ss-yuan.github.io/LMPOcc/>.

<sup>1</sup>College of Intelligent Robotics and Advanced Manufacturing, Fudan University, Shanghai, China. {ssyuan23}@m.fudan.edu.cn, {ganzhongxue, dingwenchao}@fudan.edu.cn

<sup>2</sup>State Key Laboratory of Intelligent Vehicle Safety Technology, Chongqing Changan Automobile Co., Ltd., Chongqing, China.

\*Corresponding authors: Wenchao Ding

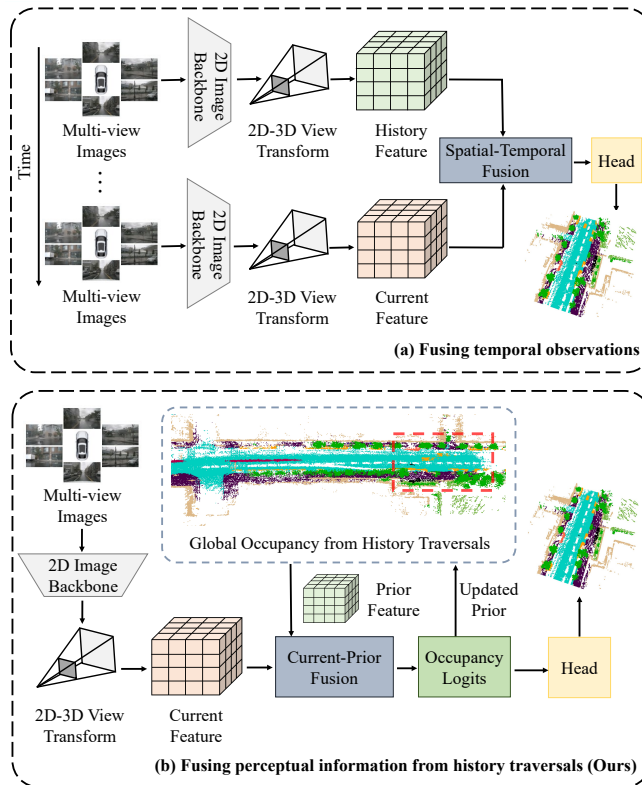


Fig. 1. Comparison of temporal information integration methods in 3D occupancy prediction. (a) Existing works primarily integrate information from adjacent observations. (b) Our work fuses perceptual information obtained from historical traversals of the current location. The historical perceptual information constructs global occupancy and serves as long-term memory priors.

alignment, self-attention mechanisms, and 3D convolution-based fusion, which primarily integrate features from adjacent frames, as illustrated in Fig. 1(a). While effective in some situations, these methods remain vulnerable when consecutive observations share similar adverse conditions such as severe occlusion or poor lighting, leading to degraded performance.

To overcome these limitations, long-term memory priors play a vital role. By aggregating perceptual information collected from repeated traversals at the same geographic locations under more favorable conditions, long-term memory priors provide rich historical context that complements and corrects deficient real-time observations. These priors naturally accumulate to form a global occupancy map, which serves as persistent memory across spatial and temporal dimensions and enables accumulation and refinement of

scene knowledge beyond single-instance observations, as shown in Fig. 1(b).

Beyond enhancing local perception, the global occupancy map offers several important benefits. It provides finer and denser geometric details than raw LiDAR data, delivering high-quality dense depth information essential for many applications. A typical example is the construction of 3D open-vocabulary maps, which support flexible and scalable scene understanding. Foundation models have achieved great success mainly on 2D images. However, to effectively leverage them in large-scale outdoor 3D scenes, reliable and dense depth information is necessary. The occupancy map reflects the scene geometry and thus provides accurate depth cues. Since real-world environments inevitably change over time, maintaining and updating the global map through continuous local occupancy predictions is important for adapting to scene dynamics.

In this paper, we propose Long-term Memory Prior Occupancy (LMPOcc), a novel framework that jointly performs local occupancy perception and constructs a global occupancy map. Specifically, our framework leverages ego-to-global coordinate transformation to simultaneously construct global occupancy representations from localized perceptual outputs, and utilize corresponding spatial memory priors from the global map to refine real-time local inference through cross-temporal feature alignment. We store visibility-region occupancy logits from each local prediction into the global map. This model-agnostic prior format enables crowdsourced construction of city-level global occupancy. To fully leverage these long-term memory priors, we design an efficient and lightweight Current-Prior Fusion module that learns adaptive weights between current and prior features to produce refined occupancy predictions. We employ ray casting to extract dense depth from occupancy maps, providing high-quality depth data for 3D open-vocabulary map construction. Extensive experiments on the Occ3D-nuScenes [1] benchmark demonstrate that LMPOcc significantly improves 3D occupancy prediction baselines and achieves state-of-the-art performance.

The main contributions of this work are summarized as follows:

- We propose LMPOcc, the first framework that leverages global occupancy as a long-term memory prior to enhance local 3D occupancy prediction while simultaneously constructing and updating the global map. Our method also provides dense depth information to support large-scale outdoor applications such as 3D open-vocabulary mapping.
- We design a plug-and-play architecture to realize the bidirectional interaction between global and local occupancy. In particular, we introduce a model-agnostic prior format and develop an efficient lightweight Current-Prior Fusion module for cross-temporal feature integration.
- Validated on the Occ3D-nuScenes benchmark, LMPOcc achieves state-of-the-art performance. We demonstrate its capability to build large-scale global occupancy maps via multi-vehicle crowdsourcing, and to leverage occupancy-derived dense depth for 3D open vocabulary map construction.

## II. RELATED WORKS

### A. 3D Semantic Occupancy Prediction

Vision-based 3D occupancy prediction has seen significant advancements driven by a variety of methodological innovations. Early supervised frameworks like MonoScene [2] establish 2D-3D U-Net architectures, while BEVDet [3] and BEVFormer [4] introduce view transformation via LSS [5] projection and spatio-temporal transformers respectively. Recent works propose strategies to enhance both representation and computational efficiency in 3D occupancy prediction. For example, TPVFormer [6] introduces three-view perspective encoding and SurroundOcc [7] employs multi-scale refinement for spatial detail capture. Unsupervised approaches like SelfOcc [8] and OccNeRF [9] circumvent dense labeling through neural rendering. To enhance computational efficiency, VoxFormer [10] and OctreeOcc [11] utilize sparse voxel representations to optimize memory and speed. GaussianFusionOcc [12] fuses multi-sensor data using semantic 3D Gaussians and deformable attention for efficient and accurate 3D occupancy prediction.

### B. Memory Fusion for 3D Perception

Current memory fusion methods for 3D perception can be divided into three paradigms, involving different ways of processing history information. Attention-based approaches [4] leverage transformers for implicit temporal modeling, demonstrating strong dependency capture while lacking explicit geometric constraints. Geometry-aligned fusion methods exemplified by BEVDet4D [13] and PanoOcc [14] employ spatial feature alignment through estimated camera poses coupled with concatenation-convolution operations, achieving computational efficiency at the expense of long-term temporal consistency. Emerging cost volume techniques address these limitations through geometric depth reasoning, as demonstrated by SOLOFusion [15] in image-space fusion and CVT-Occ [16] in 3D voxel adaptation. ST-Occ [17] proposes a scene-centered spatiotemporal memory to efficiently aggregate historical occupancy features from adjacent scenes. The methods discussed above primarily integrate information from adjacent observations, while our work leverages perceptual information acquired from historical traversals of identical geographic locations.

## III. APPROACH

### A. Overall Architecture

An overview of the Long-term Memory Prior Occupancy (LMPOcc) is presented in Fig. 2. LMPOcc extends the occupancy prediction baseline by incorporating the Long-term Memory Occupancy Priors (LMOP) module (see Sec. III-B), which strengthens local perception and facilitates the construction of global occupancy. The system receives inputs  $\mathcal{I} = \{\mathbf{I}, \mathbf{G}_{ego}\}$ , comprising surround-view images  $\mathbf{I}$  and the ego vehicle’s local-to-global coordinate transformation  $\mathbf{G}_{ego} \in \mathbb{R}^{4 \times 4}$ . The model processes surround-view images through an occupancy encoder to generate latent features. These latent features are then fed into the LMOP module to

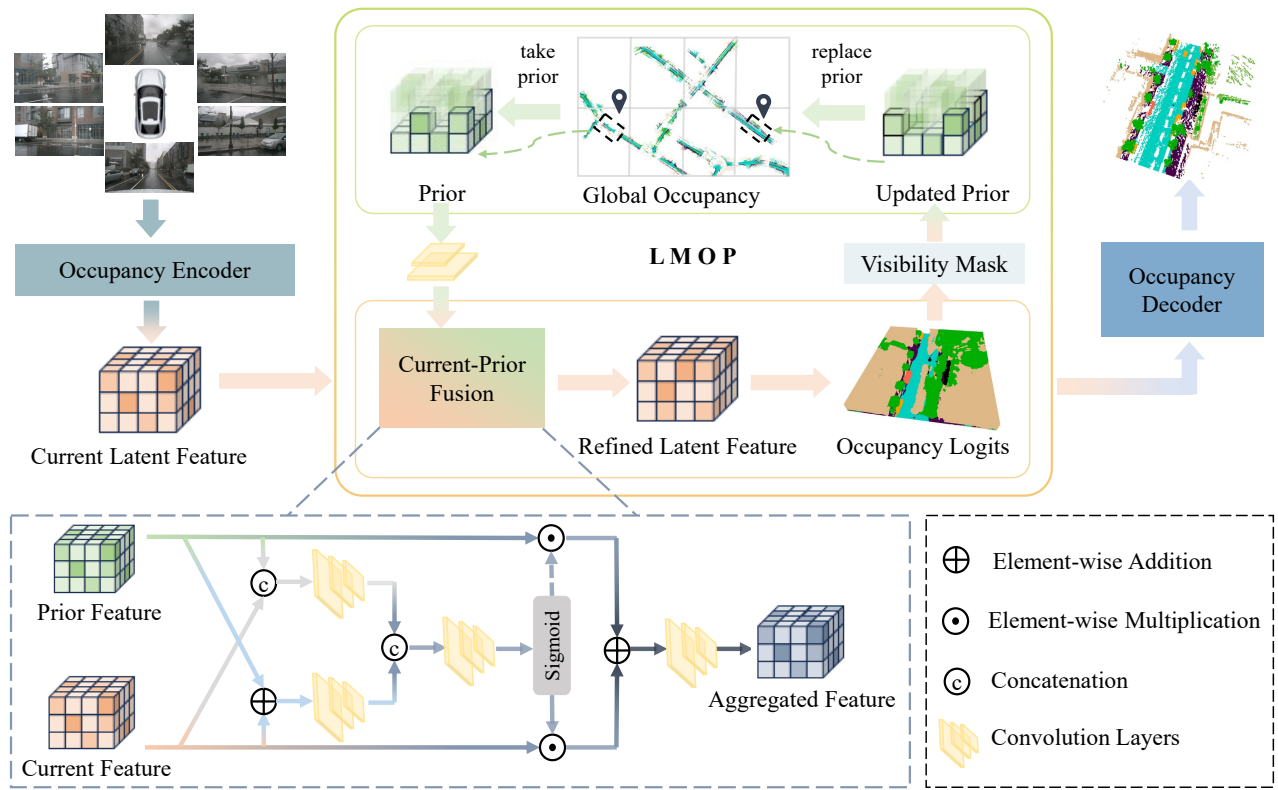


Fig. 2. An overview of our LMPOcc framework. LMPOcc firstly generates Current Latent Features from surround-view images. Then it extracts spatially-aligned Prior Features from global occupancy and integrates them via the Current-Prior Fusion Module to generate Refined Latent Features. The refined latent features decode current occupancy logits, which are stored into corresponding locations in the global occupancy after visibility masking. Existing occupancy priors at these locations are replaced by the updated logits. Finally, the occupancy logits are converted into local current occupancy prediction results.

obtain enhanced occupancy logits, which are subsequently processed by an occupancy decoder to produce the final 3D semantic occupancy prediction results. In the LMOP module, the current features are fused with corresponding prior features through the Current-Prior Fusion module, yielding the refined latent features (see Sec. III-C). These refined latent features are subsequently transformed into occupancy logits via neural network processing. The occupancy logits are utilized to update prior features and generate the final occupancy prediction results. The model-agnostic prior format is introduced in Sec. III-D. Moreover, LMPOcc also supports the construction of large-scale 3D open vocabulary maps by leveraging the dense depth information derived from occupancy results (detailed in Sec. III-E).

### B. Long-Term Memory Occupancy Priors

The Long-term Memory Occupancy Priors (LMOP) module is plug-and-play and compatible with various 3D occupancy baselines. It enables local perception and global occupancy to reinforce each other. Inspired by NMP [18], the global occupancy adopts a sparse map tile structure, where each tile corresponds to a geographically aligned patch in the global coordinate system and is initialized as empty. This sparse structure significantly reduces memory usage by only storing tiles covering navigable urban zones (e.g., roads and accessible areas), thus avoiding the need to store the entire

city map. Tiles can be efficiently retrieved and updated via coordinate-based indexing, allowing vehicles to load relevant local map areas on demand. For each tile, we maintain and iteratively update the global occupancy representation  $\mathbf{P} \in \mathbb{R}^{H_G \times W_G \times (Z \cdot N_{\text{sem}})}$ , where  $H_G$  and  $W_G$  define the spatial resolution of the city-level map tile,  $Z$  denotes vertical discretization depth, and  $N_{\text{sem}}$  corresponds to the number of distinguishable object categories. Both the global map and the local prior feature are represented in Bird's-Eye View (BEV) format through height-to-channel transformation, as expressed in FlashOcc [19]. This BEV representation not only reduces storage costs by effectively stacking height information into channels, but also enhances bidirectional local-global indexing efficiency. The local coordinate of each pixel in the BEV feature  $c_t \in \mathbb{R}^{H \times W \times 2}$  is transformed to the corresponding global coordinate  $p_t \in \mathbb{R}^{H \times W \times 2}$  through  $\mathbf{G}_{ego} \in \mathbb{R}^{4 \times 4}$ . This transformation ensures spatial alignment between locally perceived features and the global map tiles. Establishing spatial correspondence between local and global occupancy, we align prior and current feature channels via convolutional layers, and then fuse current features with prior features to enhance local perception. The enhanced perceptual outputs, represented as occupancy logits, serve as updated priors and replace the corresponding regions within the global occupancy map. This incremental update process

allows the system to refine the global occupancy over time using new local observations, thus maintaining a persistent yet dynamically evolving city-level occupancy prior.

### C. Current-Prior Fusion

Our Current-Prior Fusion (CPFusion) is shown in Fig. 2. The CPFusion module incorporates two parallel branches, comprising a concatenation branch and an element-wise addition branch. The concatenation branch concatenates the current feature  $F_c$  and prior feature  $F_p$  to form a combined feature  $F_{cat}$ , as shown in Eq. 1. Concurrently, the element-wise addition branch gets their element-wise sum results  $F_{add}$ , as shown in Eq. 2. These two features,  $F_{cat}$  and  $F_{add}$ , are then concatenated and passed through a convolution layer followed by a sigmoid activation function, yielding a tensor  $\alpha$  with values constrained between 0 and 1, as shown in Eq. 3. This tensor  $\alpha$  serves as a weighting factor to dynamically balance the contributions of the current and prior features through a weighted summation, as expressed in Eq. 4.

$$F_{cat} = W_1(f_{cat}(F_c, F_p)), \quad (1)$$

$$F_{add} = W_2(F_c + F_p), \quad (2)$$

$$\alpha = \sigma(W_3(f_{cat}(F_{cat}, F_{add}))), \quad (3)$$

$$F_{agg} = \alpha \odot F_c + (1 - \alpha) \odot F_p, \quad (4)$$

where  $f_{cat}(\cdot)$  and  $\sigma(\cdot)$  refer to the concatenation and sigmoid, respectively.  $W_1$ ,  $W_2$  and  $W_3$  denote convolution layers.  $F_{agg}$  represents the output features of CPFusion.

### D. Model-Agnostic Prior Format

The prior is stored as occupancy logits, ensuring that the global occupancy prior remains agnostic to any specific occupancy prediction model during deployment. After fusing current and prior features into refined latent features, the network outputs occupancy logits  $O_L \in \mathbb{R}^{H_L \times W_L \times Z \times N_{sem}}$ , where  $H_L$ ,  $W_L$ ,  $Z$  denote spatial dimensions and  $N_{sem}$  denotes the number of semantic classes. To avoid storing noise outside the visible regions into the prior, we employ the camera visibility mask to retain only the content within the observable regions of the occupancy logits. The camera visibility mask is generated by casting rays from each camera origin towards voxel centers, following the approach in Occ3D-NuScenes [1]. Along each ray, the first intersected occupied voxel is labeled as ‘‘observed’’, while subsequent voxels along the same ray are marked as ‘‘unobserved’’. Any voxel not intersected by these rays is automatically assigned an ‘‘unobserved’’ status. The masked occupancy logits are then reshaped to  $H_L \times W_L \times (Z \cdot N_{sem})$  and used to update the corresponding regions in the global occupancy map. Contrary to common practice, our experiments show that filtering dynamic objects from the prior fails to improve model performance, as shown in Table VI. This suggests historical dynamic objects may provide effective information for local perception. Therefore, our method retains dynamic components within the prior. We further propose two strategies for dynamic object removal and discuss the impact of dynamic components within the prior in Sec. IV-D.

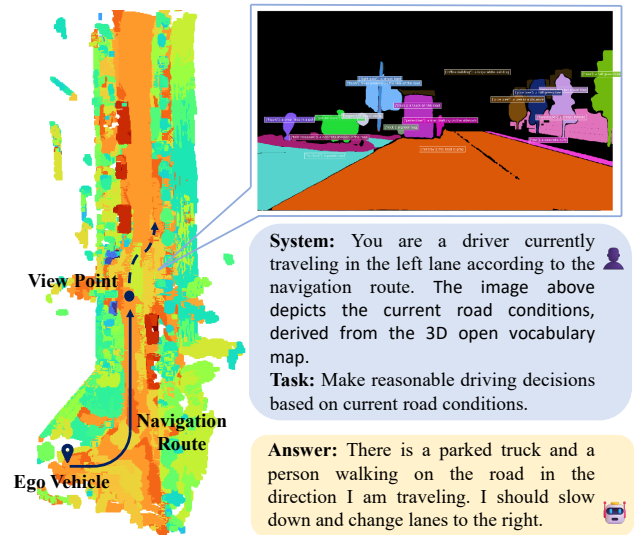


Fig. 3. Demonstration of using a 3D open vocabulary map to interact with a Vision-Language Model (VLM). A viewpoint along the navigation route is selected, from which a semantic map is rendered by the 3D open vocabulary map and then interpreted by the VLM. This enables the VLM to anticipate upcoming road conditions and make informed driving decisions to handle challenging scenarios. The example dialogue shows the VLM analyzing the scene and recommending appropriate actions.

### E. Construction of 3D Open Vocabulary Maps

To build 3D open vocabulary maps for large-scale outdoor scenes, three key input modalities need to be available, namely images, their associated poses, and corresponding depth information. Dense depth estimation is crucial for accurately projecting 2D open vocabulary semantic information into the 3D space. LMPOcc produces dense voxel occupancy grids from multi-view images, which serve as the basis for high-quality outdoor depth estimation via ray casting.

Formally, given an image pixel  $\mathbf{u} = (u, v)^\top$ , its back-projected ray direction  $\mathbf{r}_c$  in the camera coordinate system is:

$$\mathbf{r}_c = \mathbf{K}^{-1}[u, v, 1]^\top = [x_c, y_c, 1]^\top \quad (5)$$

where  $\mathbf{K}$  is the camera intrinsic matrix. Given sampled depths  $d_i = i\Delta d$ ,  $i = 1, \dots, N_d$ , 3D points along the ray in camera coordinates are computed as:

$$\mathbf{p}_{c,i} = d_i \mathbf{r}_c = [d_i x_c, d_i y_c, d_i]^\top \quad (6)$$

These points are transformed to ego coordinates in homogeneous form as:

$$\tilde{\mathbf{p}}_{ego,i} = \mathbf{T}_{camera \rightarrow ego}[\mathbf{p}_{c,i}, 1]^\top \in \mathbb{R}^4, \quad (7)$$

where  $\mathbf{T}_{camera \rightarrow ego} \in SE(3)$  is the extrinsic transformation. The corresponding Euclidean coordinates  $\mathbf{p}_{ego,i} \in \mathbb{R}^3$  are obtained by extracting the first three components of  $\tilde{\mathbf{p}}_{ego,i}$ . Indexing the occupancy voxel grid, voxel indices are computed by:

$$\mathbf{v}_i = \left\lfloor \frac{\mathbf{p}_{ego,i} - \mathbf{p}_{min}}{v_{size}} \right\rfloor, \quad (8)$$

where  $\mathbf{p}_{min}$  denotes the minimum grid coordinate, and  $v_{size}$  denotes the voxel size. Depth at pixel  $\mathbf{u}$  is defined as the

TABLE I

3D OCCUPANCY PREDICTION PERFORMANCE ON THE OCC3D-NUSCENES VALIDATION SET. BOTH THE SMALL VERSION AND LARGE VERSION OF LMPOCC OUTPERFORM THE MODELS THAT HAVE SIMILAR SETTINGS.

Method	History Frame	Resolution	Backbone	mIoU $\uparrow$	others $\uparrow$	barrier $\uparrow$	bicycle $\uparrow$	bus $\uparrow$	car $\uparrow$	cons. veh. $\uparrow$	motorcycle $\uparrow$	pedestrian $\uparrow$	traffic cone $\uparrow$	trailer $\uparrow$	truck $\uparrow$	drive. surf. $\uparrow$	other flat $\uparrow$	sidewalk $\uparrow$	terrain $\uparrow$	manmade $\uparrow$	vegetation $\uparrow$
					■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■
TPVFormer [6]	X	928 × 1600	R101	27.83	7.22	38.90	13.67	40.78	45.90	17.23	19.99	18.85	14.30	26.69	34.17	55.65	35.47	37.55	30.70	19.40	16.78
CTF-Occ [1]	X	928 × 1600	R101	28.53	8.09	39.33	20.56	38.29	42.24	16.93	24.52	22.72	21.05	22.98	31.11	53.33	33.84	37.98	33.23	20.79	18.00
OccFormer [20]	X	256 × 704	R50	20.40	6.62	32.57	13.13	20.37	37.12	5.04	14.02	21.01	16.96	9.34	20.64	40.89	27.02	27.43	18.65	18.78	16.90
FlashOcc (M1) [19]	X	256 × 704	R50	32.08	6.74	37.65	10.26	39.55	44.36	14.88	13.4	15.79	15.38	27.44	31.73	78.82	37.98	48.7	52.5	37.89	32.24
DHD-S [21]	X	256 × 704	R50	36.50	10.59	43.21	23.02	40.61	47.31	21.68	23.25	23.85	23.40	31.75	34.15	80.16	41.30	49.95	54.07	38.73	33.51
LightOcc-S [22]	X	256 × 704	R50	37.93	<b>11.72</b>	45.61	<b>25.40</b>	43.10	48.66	21.38	<b>25.58</b>	26.58	<b>29.19</b>	33.18	35.09	79.97	41.81	50.35	53.88	39.40	33.97
LMPOcc-S (Ours)	X	256 × 704	R50	<b>40.38</b>	10.97	<b>48.87</b>	23.66	<b>43.31</b>	<b>51.27</b>	<b>23.61</b>	23.79	<b>27.49</b>	26.28	<b>36.26</b>	<b>37.95</b>	<b>81.97</b>	<b>42.06</b>	<b>52.09</b>	<b>58.43</b>	<b>52.96</b>	<b>45.45</b>
FastOcc [23]	16	640 × 1600	R101	39.21	12.06	43.53	28.04	44.80	52.16	22.96	29.14	29.68	26.98	30.81	38.44	82.04	41.93	51.92	53.71	41.04	35.49
PanoOcc [14]	3	512 × 1408	R101	42.13	11.67	50.48	29.64	49.44	55.52	23.29	33.26	30.55	30.99	34.43	42.57	83.31	44.23	54.40	56.04	45.94	40.40
FB-Occ [24]	4	512 × 1408	R101	43.41	12.10	50.23	32.31	48.55	52.89	31.20	31.25	30.78	32.33	37.06	40.22	83.34	<b>49.27</b>	57.13	59.88	47.67	41.76
OctreeOcc [11]	4	512 × 1408	R101	44.02	11.96	51.70	29.93	53.52	56.77	30.83	33.17	30.65	29.99	37.76	43.87	83.17	44.52	55.45	58.86	49.52	46.33
GEOcc [25]	8	512 × 1408	SwinB	44.67	14.02	51.40	33.08	52.08	56.72	30.04	33.54	32.34	35.83	39.34	44.18	83.49	46.77	55.72	58.94	48.85	43.00
COTR [26]	8	512 × 1408	SwinB	46.20	<b>14.85</b>	53.25	<b>35.19</b>	50.83	57.25	<b>35.36</b>	34.06	33.54	<b>37.14</b>	38.99	44.97	84.46	48.73	57.60	61.08	51.61	46.72
FlashOcc (M3) [19]	1	512 × 1408	SwinB	43.52	13.42	51.07	27.68	51.57	56.22	27.27	29.98	29.93	29.80	37.77	43.52	83.81	46.55	56.15	59.56	50.84	44.67
DHD-L [21]	1	512 × 1408	SwinB	45.53	14.08	53.12	32.39	52.44	57.35	30.83	35.24	33.01	33.43	37.90	45.34	84.61	47.96	57.39	60.32	52.27	46.24
LightOcc-L [22]	1	512 × 1408	SwinB	46.00	14.50	52.27	34.45	<b>53.79</b>	57.33	31.80	<b>35.83</b>	33.60	36.09	<b>39.89</b>	46.09	84.23	48.10	57.14	60.02	51.70	45.23
LMPOcc-L (Ours)	1	512 × 1408	SwinB	<b>46.61</b>	13.68	<b>53.88</b>	31.65	53.19	<b>58.53</b>	30.68	34.7	<b>34.86</b>	34.6	39.45	<b>47.17</b>	<b>85.08</b>	47.85	<b>58.11</b>	<b>61.56</b>	<b>57.36</b>	<b>49.97</b>

smallest sampled depth  $d_i$  along the corresponding ray for which the voxel occupancy label is not free space:

$$D(u, v) = \min\{d_i \mid O(\mathbf{v}_i) \neq l_{\text{free}}\}, \quad (9)$$

where  $O(\cdot)$  returns voxel occupancy label and  $l_{\text{free}}$  indicates free space. In case no occupied voxel is found along the ray, the depth is set to a predefined maximum depth  $D_{\text{max}}$ .

Using dense depth maps and known poses, 3D open vocabulary maps can be built with existing frameworks such as OpenGraph [27]. Fig. 3 illustrates an empirical example of how 3D open vocabulary maps enable interaction with vision-language models (VLMs) and large language models (LLMs).

## IV. EXPERIMENT

### A. Datasets and Metrics

We evaluate our approach on the Occ3D-nuScenes benchmark [1], which extends the widely adopted large-scale autonomous driving dataset nuScenes [28]. This benchmark comprises 700 training scenes and 150 validation scenes, each with 40 annotated samples captured at 2Hz over 20-second sequences. The dataset spans a spatial domain of  $X \in [-40m, 40m]$  and  $Y \in [-40m, 40m]$  for horizontal dimensions, with vertical coverage from  $Z \in [-1m, 5.4m]$ . The occupancy annotations in this dataset are represented as axis-aligned voxels with 0.4m edge length, achieving a resolution of  $200 \times 200 \times 16$  voxels. Each voxel is labeled with one of 17 semantic categories or marked as free space (non-occupied). In our work, *Dynamic* represents dynamic semantic categories, which encompass *others*, *barrier*, *bicycle*, *bus*, *car*, *construction vehicle*, *motorcycle*, *pedestrian*, *traffic cone*, *trailer* and *truck*.

*Static* denotes static semantic categories, which comprise *driveable surface*, *other flat*, *sidewalk*, *terrain*, *manmade* and *vegetation*. For performance evaluation, we employ the mean Intersection-Over-Union (mIoU) metric aggregated across all semantic classes.

### B. Implementation Details

We employ FlashOcc [19] and DHD [21] as baseline models. During training, we initialize the models with their pretrained weights and freeze the parameters preceding the current latent feature. After integrating our LMOP module, we train for another 24 epochs while maintaining identical experimental configurations to the baseline setup. We disable BEV-space data augmentation to prevent misalignment between current features and prior features. The channel dimension of the global occupancy in LMPOcc is computed as the height of the occupancy multiplied by the number of semantic categories, specifically  $16 \times 18$ , while other configurations regarding the global map remain consistent with those in Neural Map Prior [18]. All models are trained with a batch size of 4 on 6 NVIDIA A100 GPUs. The 3D open vocabulary maps are constructed using 6-view images per occupancy frame, with depth sampled at intervals  $\Delta d = 0.1m$  up to a maximum depth  $D_{\text{max}} = 100m$ .

### C. Main Results

Our Long-term Memory Occupancy Prior (LMOP) serves as a plug-and-play approach applicable to diverse occupancy algorithms. To illustrate this, we integrate LMOP into two base models: FlashOcc and DHD. We use the same hyperparameter settings as in their original designs. During training, we freeze all the modules preceding the current

TABLE II

THE PERFORMANCE OF OCCUPANCY PREDICTION METHODS AND THEIR LMOP VERSIONS ON THE OCC3D-NUSCENES VALIDATION SET. BY ADDING LONG-TERM MEMORY KNOWLEDGE, LMOP CONSISTENTLY IMPROVES THESE METHODS.

Model	mIoU		
	Dynamic	Static	All
FlashOcc-M0	23.67	47.11	31.94
FlashOcc-M0 + LMOP	<b>26.36</b>	<b>53.29</b>	<b>35.87</b>
$\Delta$ mIoU	+2.69	+6.18	+3.93
DHD-S	29.35	49.62	36.5
DHD-S + LMOP	<b>32.13</b>	<b>55.49</b>	<b>40.38</b>
$\Delta$ mIoU	+2.78	+5.87	+3.88

latent features and only train the downstream modules, especially LMOP. As evidenced in Table II, LMOP consistently improves occupancy prediction compared to baseline models.

We compare LMPOcc with the state-of-the-art occupancy prediction method on the Occ3D-nuScenes benchmark [1]. The experiment results in Table I show that both LMPOcc-S and LMPOcc-L outperform other approaches that have similar model settings. The baseline of LMPOcc-S and LMPOcc-L are respectively DHD-S and DHD-L. These persuasive experiment results highlight the effectiveness of LMOP in occupancy prediction, especially on static semantic categories.

#### D. Ablation Studies

For efficient validation, we conduct ablation studies on LMPOcc-S with DHD-S as the baseline. The ablation studies are conducted on Current-Prior Fusion, Visibility Mask, and Dynamic Removal.

**Current-Prior Fusion.** The performance of different fusion methods is presented in Table III. The proposed Current-Prior Fusion module demonstrates superior performance compared to both direct concatenation and element-wise addition of current and prior features. Moreover, the joint utilization of both the Concatenation Branch and Addition Branch within the Current-Prior Fusion module yields the optimal performance. In Table III, ‘C.B.’ and ‘A.B.’ represent the Concatenation Branch and Addition Branch within our Current-Prior Fusion module.

Neural Map Prior [18] employs cross-attention coupled with Gated Recurrent Unit (GRU) to fuse current features with prior features. As illustrated in Table IV, our Current-Prior Fusion method not only surpasses Neural Map Prior’s fusion modules in performance, but also reduces computational latency, thus demonstrating significant practical advantages.

**Visibility Mask.** We analyze the impact of the visibility mask within our framework. As demonstrated in Table V, LMPOcc slightly underperforms the baseline without the visibility mask, whereas its performance significantly surpasses the baseline when the visibility mask is applied. This is because regions outside the visibility mask contain

TABLE III

ABLATION STUDY OF CURRENT-PRIOR FUSION. THE TERM ‘CONCAT’ DENOTES DIRECT CONCATENATION OF CURRENT AND PRIOR FEATURES. ‘ADD’ INDICATES ELEMENT-WISE ADDITION OF THE TWO FEATURES. ‘C.B.’ AND ‘A.B.’ REPRESENT THE CONCATENATION BRANCH AND ADDITION BRANCH WITHIN OUR CURRENT-PRIOR FUSION MODULE.

Method		mIoU		
		Dynamic	Static	All
Concat		31.51	54.26	39.54
Add		31.52	54.41	39.60
CPFusion	C.B.			
	✓	31.95	55.20	40.16
	✓	31.76	55.09	40.00
	✓	<b>32.13</b>	<b>55.49</b>	<b>40.38</b>

TABLE IV

COMPARATIVE ANALYSIS OF FUSION MODULES IN NEURAL MAP PRIOR. LATENCIES ARE EVALUATED ON A SINGLE A100 GPU.

Method	mIoU $\uparrow$	Latency $\downarrow$
Cross Attention + GRU [18]	39.80	11.6 ms
CPFusion (Ours)	<b>40.38</b>	<b>7.1 ms</b>

TABLE V

ABLATION STUDY OF CAMERA VISIBILITY MASK.

Method	mIoU		
	Dynamic	Static	All
Baseline	29.35	49.62	36.5
Baseline + LMOP (w/o Mask)	28.87	49.43	36.13
Baseline + LMOP (w/ Mask)	<b>32.13</b>	<b>55.49</b>	<b>40.38</b>

TABLE VI

ANALYSIS OF DYNAMIC TARGETS IN LMOP.

Method	mIoU		
	Dynamic	Static	All
Removing Dynamic v1	30.45	55.28	39.21
Removing Dynamic v2	32.05	55.29	40.25
Retaining Dynamic	<b>32.13</b>	<b>55.49</b>	<b>40.38</b>

substantial noise, thus storing information exclusively within the visible regions ensures the validity of the prior.

**Discussion of Dynamic Targets.** Given the temporal inconsistencies of dynamic targets at the same geographic location across distinct observation periods, we conduct analysis on the handling of dynamic targets within the global prior in our framework. We propose two distinct methodologies for eliminating dynamic elements within the prior framework. After applying the argmax operation to the occupancy logits, the semantic label for each voxel can be obtained. We create free masks and dynamic masks to, respectively, mask out the free components and dynamic components from the occupancy logits. The dynamic components comprise dynamic semantic categories, as detailed in Sec. IV-A. The first method sets the dynamic components in the occupancy logits as zero, denoted as *Removing Dynamic v1* in Table VI.

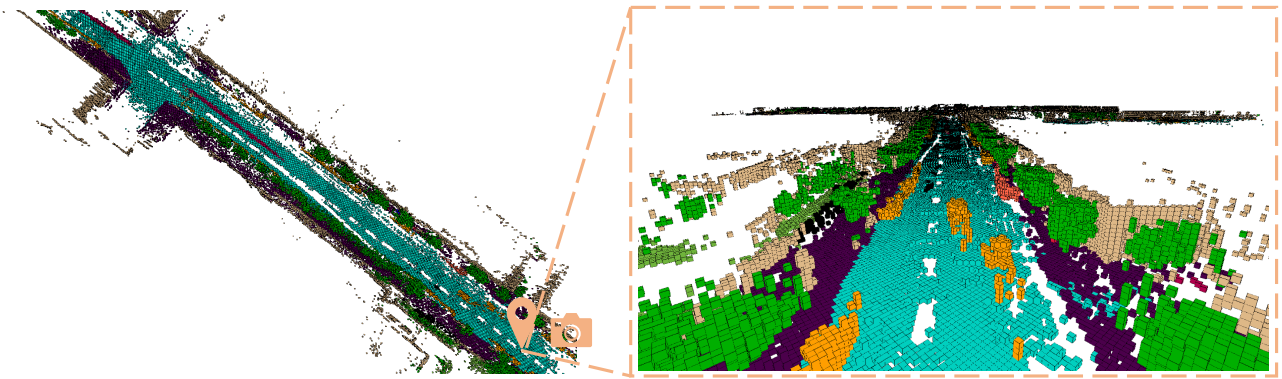


Fig. 4. Visualization results of a region within our global occupancy. The left side shows the top view, and the right side shows the front view.

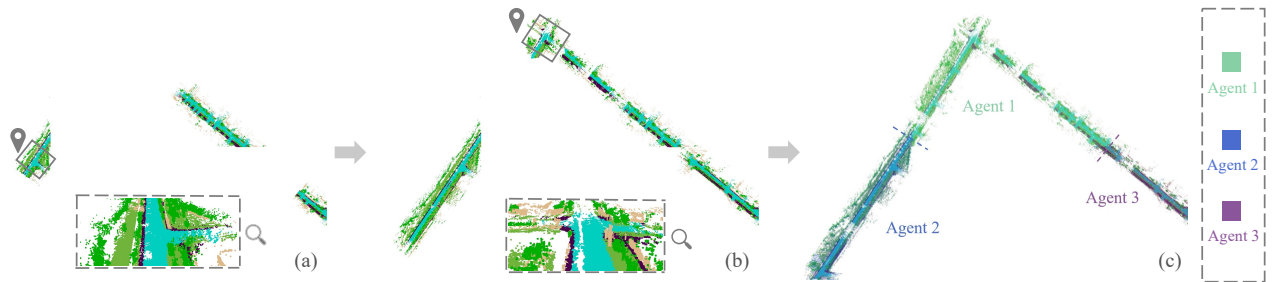


Fig. 5. Visualization results of global occupancy construction via crowdsourcing methodologies. Three collaborative agents construct the global occupancy map through crowdsourcing. (a) and (b) show the intermediate stages of the occupancy construction process. (c) displays the crowdsourced mapping result. Three colors mark the areas mapped by each agent.

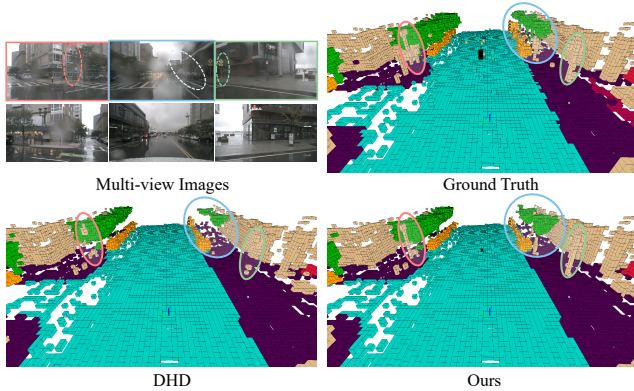


Fig. 6. Visualization results of LMPOcc and DHD. It is a low-visibility rainy scene. LMPOcc leverages long-term memory priors to detect objects not visible in current sensory observations, demonstrating significant improvement over the baseline. The color-semantic category correspondence is detailed in Table I.

The second method randomly selects the free components in the occupancy logits and replaces the dynamic components with them, denoted as *Removing Dynamic v2* in Table VI.

The results indicate that removing dynamic elements does not yield performance improvements, as detailed in Table VI. A plausible explanation lies in our Current-Prior Fusion module’s capability to adaptively process dynamic components within prior features. Another critical factor is that dynamic objects typically exhibit spatial-temporal distribution patterns in specific regions, which serves as effective prior knowledge

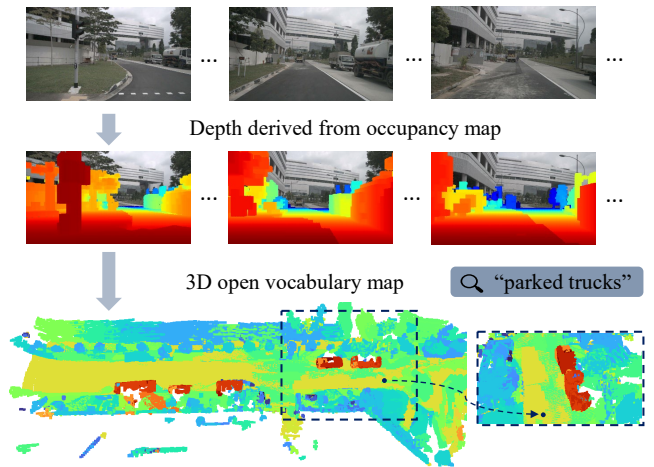


Fig. 7. Visualization results of the 3D open vocabulary map. Depth information is derived from the occupancy map, and the outdoor open-vocabulary map is constructed based on the OpenGraph method. The bottom row shows scene-level occupancy grounding for specific queries (e.g., “parked trucks”).

to enhance dynamic object detection. This mirrors how the human brain leverages spatial priors in perceptual processing to identify moving entities.

### E. Visualization

Our global occupancy is shown in Fig. 4. While enhancing local perception capabilities, LMPOcc can construct global

occupancy for large-scale scenes. As shown in Fig. 5, three collaborative agents construct the global occupancy map through crowdsourcing. The visualization results of LMPOcc-S are presented in Fig. 6. It shows a rainy, low-visibility scene where LMPOcc leverages long-term memory priors to detect objects missing in current views, significantly outperforming the baseline.

#### F. Construction of 3D Open Vocabulary Maps

In Fig. 7, we showcase the results of a 3D open-vocabulary map constructed with OpenGraph [27], utilizing depth generated from the occupancy map for a sequence of multi-view images. The first row shows the image sequence used for building the 3D map, the second row illustrates the depth extracted from the occupancy map for each image via ray casting, and the third row displays a scene-level 3D open-vocabulary map along with the grounding of "parked trucks" on the map.

### V. CONCLUSION

We propose LMPOcc, a novel 3D occupancy prediction framework that leverages long-term memory priors from historical traversals to enhance local perception and build unified global occupancy maps. By introducing a model-agnostic prior format, LMPOcc ensures compatibility across different prediction baselines, while a lightweight Current-Prior Fusion module adaptively integrates prior and current features. Validated on the Occ3D-nuScenes benchmark, LMPOcc achieves state-of-the-art local occupancy prediction, especially for static semantic classes. Additionally, it supports large-scale global occupancy construction via multi-vehicle crowdsourcing and facilitates 3D open-vocabulary map building through occupancy-derived dense depth. Future research could conduct in-depth investigations into dynamic object processing within prior information and develop learnable optimization frameworks for global occupancy.

### REFERENCES

- [1] X. Tian, T. Jiang, L. Yun, Y. Mao, H. Yang, Y. Wang, Y. Wang, and H. Zhao, "Occ3d: A large-scale 3d occupancy prediction benchmark for autonomous driving," *Advances in Neural Information Processing Systems*, vol. 36, pp. 64 318–64 330, 2023.
- [2] A.-Q. Cao and R. De Charette, "Monoscene: Monocular 3d semantic scene completion," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 3991–4001.
- [3] J. Huang, G. Huang, Z. Zhu, Y. Ye, and D. Du, "Bevdet: High-performance multi-camera 3d object detection in bird-eye-view," *arXiv preprint arXiv:2112.11790*, 2021.
- [4] Z. Li, W. Wang, H. Li, E. Xie, C. Sima, T. Lu, Q. Yu, and J. Dai, "Bevformer: learning bird's-eye-view representation from lidar-camera via spatiotemporal transformers," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- [5] J. Phillion and S. Fidler, "Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d," in *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16*. Springer, 2020, pp. 194–210.
- [6] Y. Huang, W. Zheng, Y. Zhang, J. Zhou, and J. Lu, "Tri-perspective view for vision-based 3d semantic occupancy prediction," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2023, pp. 9223–9232.
- [7] Y. Wei, L. Zhao, W. Zheng, Z. Zhu, J. Zhou, and J. Lu, "Surroundocc: Multi-camera 3d occupancy prediction for autonomous driving," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 21 729–21 740.

- [8] Y. Huang, W. Zheng, B. Zhang, J. Zhou, and J. Lu, "Selfocc: Self-supervised vision-based 3d occupancy prediction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 19 946–19 956.
- [9] C. Zhang, J. Yan, Y. Wei, J. Li, L. Liu, Y. Tang, Y. Duan, and J. Lu, "Occnerf: Self-supervised multi-camera occupancy prediction with neural radiance fields," *CoRR*, 2023.
- [10] Y. Li, Z. Yu, C. Choy, C. Xiao, J. M. Alvarez, S. Fidler, C. Feng, and A. 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*, 2023, pp. 9087–9098.
- [11] Y. Lu, X. Zhu, T. Wang, and Y. Ma, "Octreeocc: Efficient and multi-granularity occupancy prediction using octree queries," *arXiv preprint arXiv:2312.03774*, 2023.
- [12] T. Pavković, M.-A. N. Mahani, J. Niedermayer, and J. Betz, "Gaussianfusionocc: A seamless sensor fusion approach for 3d occupancy prediction using 3d gaussians," *arXiv preprint arXiv:2507.18522*, 2025.
- [13] J. Huang and G. Huang, "Bevdet4d: Exploit temporal cues in multi-camera 3d object detection," *arXiv preprint arXiv:2203.17054*, 2022.
- [14] Y. Wang, Y. Chen, X. Liao, L. Fan, and Z. Zhang, "Panoocc: Unified occupancy representation for camera-based 3d panoptic segmentation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2024, pp. 17 158–17 168.
- [15] J. Park, C. Xu, S. Yang, K. Keutzer, K. Kitani, M. Tomizuka, and W. Zhan, "Time will tell: New outlooks and a baseline for temporal multi-view 3d object detection," *arXiv preprint arXiv:2210.02443*, 2022.
- [16] Z. Ye, T. Jiang, C. Xu, Y. Li, and H. Zhao, "Cvt-occ: Cost volume temporal fusion for 3d occupancy prediction," in *European Conference on Computer Vision*. Springer, 2024, pp. 381–397.
- [17] Z. Leng, J. Yang, W. Yi, and B. Zhou, "Occupancy learning with spatiotemporal memory," *arXiv preprint arXiv:2508.04705*, 2025.
- [18] X. Xiong, Y. Liu, T. Yuan, Y. Wang, Y. Wang, and H. Zhao, "Neural map prior for autonomous driving," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 17 535–17 544.
- [19] Z. Yu, C. Shu, J. Deng, K. Lu, Z. Liu, J. Yu, D. Yang, H. Li, and Y. Chen, "Flashocc: Fast and memory-efficient occupancy prediction via channel-to-height plugin," *arXiv preprint arXiv:2311.12058*, 2023.
- [20] Y. Zhang, Z. Zhu, and D. Du, "Occformer: Dual-path transformer for vision-based 3d semantic occupancy prediction," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 9433–9443.
- [21] Y. Wu, Z. Yan, Z. Wang, X. Li, L. Hui, and J. Yang, "Deep height decoupling for precise vision-based 3d occupancy prediction," *arXiv preprint arXiv:2409.07972*, 2024.
- [22] J. Zhang, Y. Zhang, Q. Liu, and Y. Wang, "Lightweight spatial embedding for vision-based 3d occupancy prediction," *arXiv preprint arXiv:2412.05976*, 2024.
- [23] J. Hou, X. Li, W. Guan, G. Zhang, D. Feng, Y. Du, X. Xue, and J. Pu, "Fastocc: Accelerating 3d occupancy prediction by fusing the 2d bird's-eye view and perspective view," *arXiv preprint arXiv:2403.02710*, 2024.
- [24] Z. Li, Z. Yu, D. Austin, M. Fang, S. Lan, J. Kautz, and J. M. Alvarez, "Fb-occ: 3d occupancy prediction based on forward-backward view transformation," *arXiv preprint arXiv:2307.01492*, 2023.
- [25] X. Tan, W. Wu, Z. Zhang, C. Fan, Y. Peng, Z. Zhang, Y. Xie, and L. Ma, "Geocc: Geometrically enhanced 3d occupancy network with implicit-explicit depth fusion and contextual self-supervision," *arXiv preprint arXiv:2405.10591*, 2024.
- [26] Q. Ma, X. Tan, Y. Qu, L. Ma, Z. Zhang, and Y. Xie, "Cotr: Compact occupancy transformer for vision-based 3d occupancy prediction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 19 936–19 945.
- [27] Y. Deng, J. Wang, J. Zhao, X. Tian, G. Chen, Y. Yang, and Y. Yue, "Opengraph: Open-vocabulary hierarchical 3d graph representation in large-scale outdoor environments," *IEEE Robotics and Automation Letters*, 2024.
- [28] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nusenes: A multimodal dataset for autonomous driving," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 11 621–11 631.